

# NCBI Molecular Biology Resources

USDA  
AFRS

*April 24, 2002*

NCBI

## NCBI Resources

- About NCBI
- NCBI Sequence Databases
  - Primary Database - GenBank
  - Derivative Databases - RefSeq
- Entrez Databases and Text Searching
- BLAST Services
- Genomic Resources

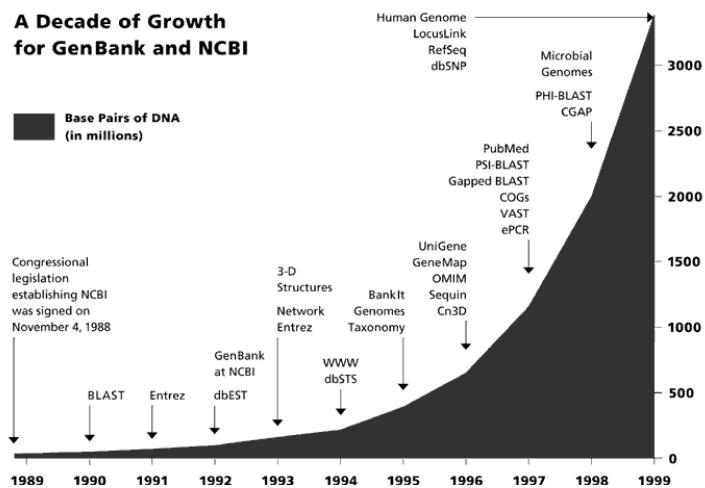
NCBI

## The National Center for Biotechnology Information (NCBI)

- Created as a part of NLM in 1988
  - Establish public databases
  - Research in computational biology
  - Develop software tools for sequence analysis
  - Disseminate biomedical information
- Tools: BLAST(1990), Entrez (1992)
- GenBank (1992)
- Free MEDLINE (PubMed, 1997)
- Human genome (2001)

NCBI

## NCBI History



NCBI

## Molecular Databases

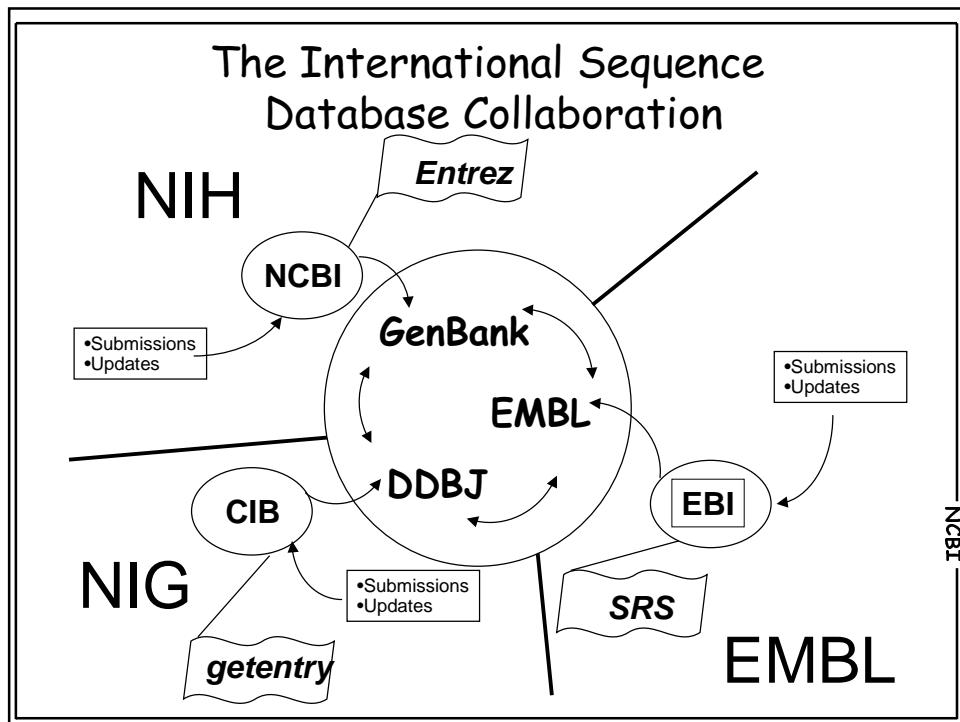
- Primary Databases
  - Original submissions by experimentalists
  - Database staff organize but don't add additional information
    - Example: GenBank
- Derivative Databases
  - Human curated
    - compilation and correction of data
    - Example: SWISS-PROT, NCBI RefSeq mRNA
  - Computationally Derived
    - Example: UniGene
  - Combinations
    - Example: NCBI Genome Assembly

NCBI

### What is GenBank? NCBI's Primary Sequence Database

- Nucleotide only sequence database
- Archival in nature
- GenBank Data
  - Direct submissions individual records (BankIt, Sequin)
  - Batch submissions via email (EST, GSS, STS)
  - ftp accounts sequencing centers
- Data shared three collaborating databases
  - GenBank
  - DNA Database of Japan (DDBJ).
  - European Molecular Biology Laboratory Database (EMBL) at EBI.

NCBI



**GenBank: NCBI's Primary Sequence Database**

RELEASE 81.0	
<b>Release 128</b>	<b>February 2002</b>
15,465,325	Records
17,089,143,893	Nucleotides
110,000 +	Species

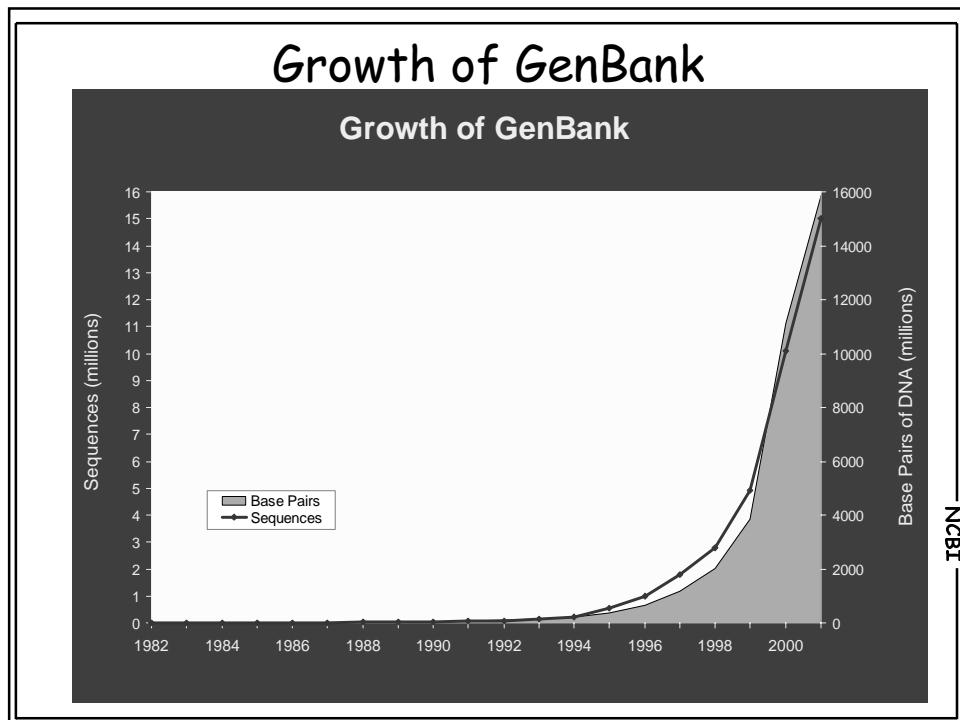
- full release every two months
- incremental and cumulative updates daily
- available only through internet

NCBI

National Library of Medicine, National Institutes of Health  
8600 Rockville Pike, Bethesda, MD 20894

ftp://ftp.ncbi.nih.gov/genbank/

60 Gigabytes of data



GenBank Divisions					
Bulk Sequence Divisions					
PAT	Patent				
EST	Expressed Sequence Tags (142 files)				
STS	Sequence Tagged Sites				
GSS	Genome Survey Sequences (48 files)				
HTG	High Throughput Genome (26 files)				
HTC	High Throughput cDNA				
CON	Contig				
Traditional Divisions					
BCT	INV	MAM	PHG	PLN	PRI
ROD	SYN	UNA	VRL	VRT	

## A Traditional GenBank Record

LOCUS	AF153828	1586 bp	mRNA	linear	PLN	18-APR-2000
LOCUS	AF153828	1586 bp	mRNA	linear	PLN	18-APR-2000
VERSION	AF153828.1	LEN:7532798				
<b>Locus Name</b>						
SOURCE		<b>length</b>				
ORGANISM						
REFERENCE			<b>molecule type</b>			
AUTHORS			mRNA==cDNA			
TITLE			DNA==gDNA			
JOU						<b>Accession Number</b>
ACCESSION	AF153828					
MEI						
VER			AF153828.1	GI:7532799		
REF						
REFERENCE	2 (bases 1..1586)					
AUTH						<b>GI Number</b>
TITLE						
JOURNAL						

## GenBank Record: Feature Table

FEATURES	Location/Qualifiers
source	1..1586 /organism="Malus x domestica" /cultivar="Granny Smith" /db_xref="taxon:3750" /tissue_type="fruit" /note="isolated from young 8 week-old fruit and fruit storage at 0.5 degrees Celsius."
/protein_id="AAC16332.2"	
/db_xref="GI:7144485"	<b>GenPept Protein IDs</b>
	/note="alpha-amylase by similarity" /codon_start=1 /product="alpha-amylase" /protein_id="AAF63239.1" /db_xref="GI:7532799" /translation="MGYGSNDSRENAQQTDIGAAVRNGREILLQAFNWESHKHDWWRN
BASE COUNT	CPAGREWTLATCGHRYAVWNK" 474 a 311 c 370 g 431 t
ORIGIN	
	1 tgcaatccgg ggccgagttg ggaaactaca tcctgagtca aatgggttac ggaagtaatg 1561 tagtgccta aaaaaaaaaaaa aaaaaaa
//	

## EST Division: Expressed Sequence Tags

```
>IMAGE:275615 5' mRNA sequence  
GACAGCATTGGGCCAGATGTCTCGCTCCGTGGCCTTAGCTGTGCTCGCCTACTCTCTCTTCTGGCC  
TGGAGGTATCCAGCGTACTCCAAAGATTCAAGTTACTCACGTACATCCAGCAGAGAATGAAAGTCAAAT  
TTCCCTGAATTGCTATGTGCTGGGTTCATCCATCCGACATTGAAGTGTACTGAAGAATGGAGAGA  
GAATTGAAAAAGTGGAGCATTCAAGCTTGCTTCAAGCAAGGACTGGTCTTCTATCTCTGTACTACAC  
TGAATTCAACCCCCACTGAAAAGATGAGTATGAGCTGGCTGGTGAACCAGTNGACTTTGTACAGNCC  
AAGTTNAGTTAACGGGNATCGAGACATGTAAGGCAGGCATATGGGAGGTTGAAGNATGCCGNTT  
TTGGATTGGATGAATTCCAATTCTGGTTGCTGNTTTAATATTGGATATGCTTTG
```

*denos*

```
>IMAGE:275615 3', mRNA sequence  
NNTCAGTTTATGATTATTTAACTTGTGGAACAAAATAACCAAGATTAACCACAACCATGCCTTACT  
TTATCAAATGTATAAGANGTAAATATGAATCTTATATGACAAAATGTTCAATTCAATTATAACAAATTCC  
AATAATCCTGTCAATNATATTCTAAATTTCACCAATTCTAAGCAGAGTATGTAATTGGAGTTAA  
CTTATGCACGCTTAACATCTAACAGCTTGAGTGCAAGAGATTGANGAGTTCAAATCTGACCAAGAT  
GTTGATGTTGGATAAGAGAATTCTGCTCCCCACCTCTANGTTGCCAGCCCTC
```

make cDNA library

80-100,000 unique cDNA clones in library

NCBI

## What is UniGene?

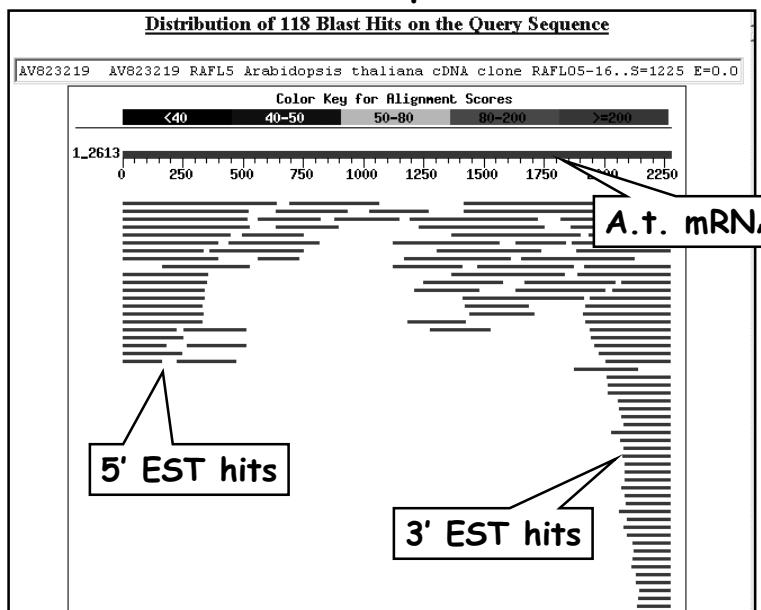
A gene-oriented view of sequence entries

- MegaBlast based automated sequence clustering
- Nonredundant set of gene oriented clusters
- Each cluster a unique gene
- Information on tissue types and map locations
- Includes well-characterized genes and novel ESTs
- Useful for gene discovery and selection of mapping reagents

NCBI

<http://www.ncbi.nlm.nih.gov/UniGene/>

## EST hits A.t. serine protease mRNA



## Arabidopsis UniGene Statistics

39,855	mRNAs + gene CDSS	UniGene Build 14 Apr. 9th, 2002		
87,006	EST, 3'reads			
42,137	EST, 5'reads			
+ 32,571	EST, other/unknown			
<hr/>				
201,569	total sequences in clusters			
 Final Number of Clusters (sets)				
<hr/>				
26,808	sets total			
	115,000,000 bp			
25,474	25,498 expected genes	: one known gene		
17,654	5% uncharacterized transcripts	: one EST		
16,326	sets contain both genes and ESTs			

## Hs UniGene Statistics

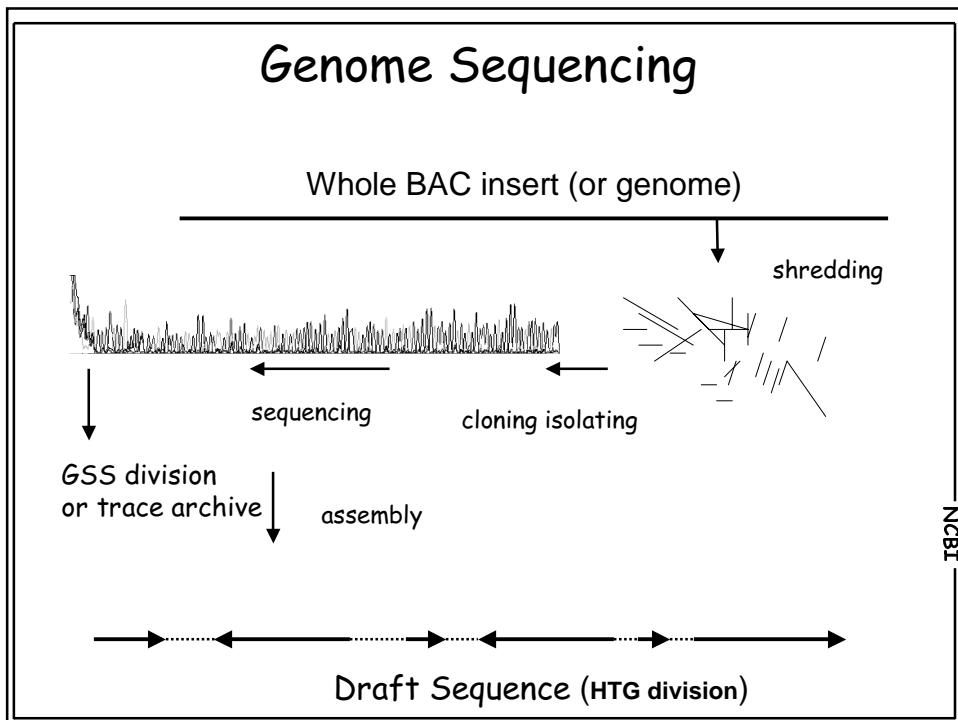
<p>73,419            mRNAs + gene CDSs          1,181,855        EST, 3' reads          1,461,928        EST, 5' reads          + 616,609        EST, other/unknown          -----          3,333,811        total sequences in clusters</p> <p><b>Final Number of Clusters (sets)</b></p> <p>=====</p> <p>98,816    sets total</p> <p>3,000,000 base pairs</p> <p>22,431    ≈ 30 K expected genes      one known gene</p> <p>97,618    ≈ 80% uncharacterized transcripts    one EST</p> <p>21,233    sets contain both genes and ESTs</p>	<div style="border: 1px solid black; padding: 2px;">           UniGene Build 148            Apr. 8th, 2002         </div>
--	---

NCBI

## UniGene Collections *Apr, 2002*

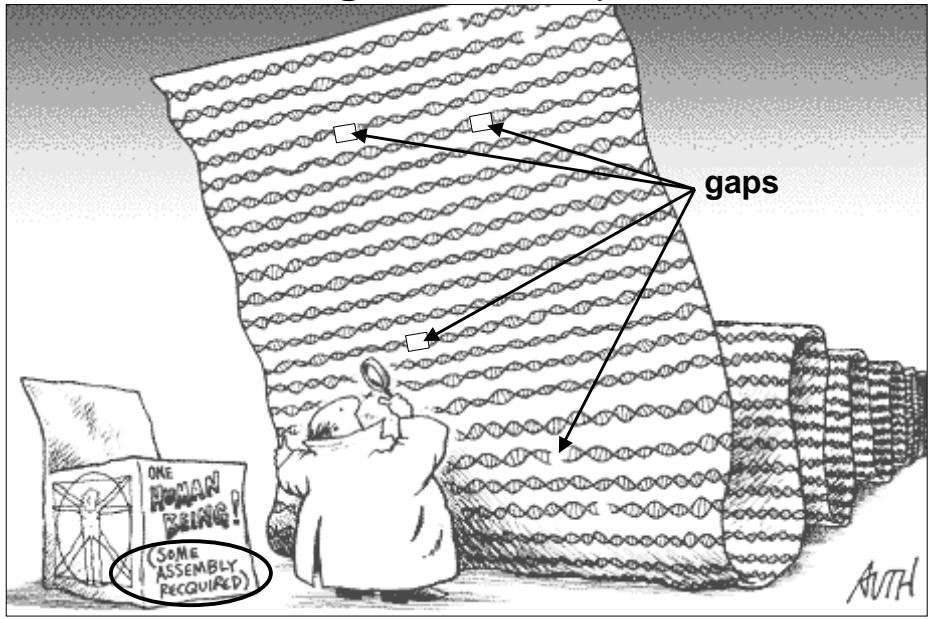
		Sequences	Clusters
<b>Animals</b>			
<i>Homo sapiens</i>	human	3,333,811	98,816
<i>Mus musculus</i>	mouse	2,274,640	86,897
<i>Rattus norvegicus</i>	rat	308,877	59,882
<i>Danio rerio</i>	zebrafish	159,261	14,893
<i>Bos taurus</i>	cow	122,503	9,303
<i>Xenopus laevis</i>	frog	120,489	16,489
<i>Anopheles gambiae</i>	mosquito	42,590	2,414
<b>Plants</b>			
<i>Arabidopsis thaliana</i>	thale cress	202,099	26,794
<i>Oryzia sativa</i>	rice	77,376	15,283
<i>Triticum aestivum</i>	wheat	35,387	3,091
<i>Hordeum vulgare</i>	barley	108,658	6,984
<i>Zea mays</i>	maize (corn)	108,030	9,889

NCBI



GSS division or trace archive		ces	
LOCUS	BH245187	195 bp	DNA linear GSS 13-NOV-2001
DEFINITION	AUIDA66TF	AUID Arabidopsis thaliana genomic clone AUIDA66, DNA sequence.	
ACCESSION	BH245187		
VERSION	BH245187.1	GI:16922701	
KEYWORDS	GSS.		
SOURCE	thale cress.		
ORGANISM	Arabidopsis thaliana		
	Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; eudicotyledons; core eudicots; Rosidae; eurosids II; Brassicales; Brassicaceae; Arabidopsis.		
REFERENCE	1 (bases 1 to 195)		
AUTHORS	Town,C.D., Whitelaw,C.A., Pai,G., Van Aken,S.E., Utterback,T.V., Feldblyum,T.V. and Fraser,C.M.		
TITLE	Survey sequencing of Arabidopsis thaliana BAC F5I22		
JOURNAL	Unpublished (2001)		
COMMENT	Contact: Chris Town TIGR 9712 Medical Center Drive, Rockville, MD 20850, USA. Tel: 301-838-3523 Fax: 301-838-0208 Email: cdtown@tigr.org From Wash. U contig 720. Seq primer: TF Class: sheared ends.		
FEATURES	Location/Qualifiers		
source	1..195 /organism="Arabidopsis thaliana" /strain="Columbia" /db_xref="taxon:3702" /clone="AUIDA66" /clone_lib="AUID" /note="Vector: pHOS2; Site_1: BstXI; 2-3 kb sheared BAC DNA inserted into pHOS2 using BstXI linkers"		

## Working Draft Sequence



## HTG Division: High Throughput Genome

phase 1 → ← → ← HTG

ACCESSION AC006228  
VERSION AC006228.1 GI:4056404  
KEYWORDS HTG; HTGS\_PHASE1.

phase 2 → → → → HTG

ACCESSION AC006228  
VERSION AC006228.2 GI:4309686  
KEYWORDS HTG; HTGS\_PHASE2.

phase 3 → → → PLN

ACCESSION AC006228  
VERSION AC006228.4 GI:4580732  
KEYWORDS HTG.

40,000 to > 350,000 bp

NCBI

## RefSeq: NCBI's Derivative Sequence Database

- Curated transcripts and proteins
  - reviewed
  - human, mouse, rat, fruit fly, zebrafish, arabidopsis
- Human model transcripts and proteins
- Assembled Genomic Regions (contigs)
  - draft human genome
  - mouse genome
- Chromosome records
  - microbial
  - organelle

NCBI

## The RefSeq Accession Numbers

### NCBI Reference Sequences

#### mRNAs and Proteins

NM\_123456 Curated mRNA }  
NP\_123456 Curated Protein }  
XM\_123456 Predicted Transcript (human)  
XP\_123456 Predicted Protein (human)

human  
mouse  
rat  
fruit fly  
zebrafish  
Arabidopsis

#### Gene Records

NG\_123456 Reference Genomic Sequence (human)

#### Assemblies

NT\_123456 Contig (Mouse and Human)

NC\_123455 Chromosome (Microbial, Arabidopsis )

NCBI

## GenBank Sequences: human CFTR

Show [10] ▾ Items 1-83 of 83 One page

1: M86631 Related Sequences, PubMed, Taxonomy  
Homo sapiens (clone ST-18-5916) cystic fibrosis transmembrane conductance regulator (CFTR) gene, 3' end intron 17B; complete exon 18; complete intron 18  
gi|80296|gb|M86631.1|HUMCFTR[180296]

2: M64699 Related Sequences, OMIM, Protein, PubMed, Taxonomy  
Homo sapiens cystic fibrosis transmembrane conductance regulator isoform 36 (CFTR) mRNA, partial cds  
gi|408281|gb|M64699.1|M64699[408281]

3: AL121762 Related Sequences, Taxonomy  
Human DNA sequence from clone RP4-410C12 on chromosomes 20 Contains the 3' end of a novel gene, a putative novel gene, a pseudogene similar to part of the cystic fibrosis transmembrane conductance regulator (CFTR), four CpG islands, ESTs, STSs and Oligos, complete sequence  
gi|8574423|emb|AL121762.1|HSJ810C12[8574423]

4: AH006034 OMIM, Protein, PubMed, Taxonomy  
Human cystic fibrosis transmembrane conductance regulator (CFTR) gene  
gi|10637|gb|AH006034.1|SEQ\_1|HUMCFTR[306337]

5: M55131 Related Sequences, ProbeSet, OMIM, Protein, PubMed, Taxonomy  
Human cystic fibrosis transmembrane conductance regulator (CFTR) gene, exon 24  
gi|10636|gb|M55131.1|HUMCFTR[120|306336]

6: M55130 Related Sequences, Protein, PubMed, Taxonomy  
Human cystic fibrosis transmembrane conductance regulator (CFTR) gene, exon 23



## Curated RefSeq Records: NM\_, NP\_

LOCUS	NM_000492	6159 bp	mRNA	PRI	26-JUL-1999
DEFINITION	Homo sapiens cystic fibrosis transmembrane conductance regulator (CFTR) mRNA				
AC	<u>REFSEQ:</u> This reference sequence was derived from M28668.1, M55131.1.				
	e				
	On Feb 17, 2000 this sequence version replaced gi:4502784.				
LO	Summary: Cystic fibrosis transmembrane conductance regulator is member 7 of the ATP-binding cassette sub-family C. The protein				
DE	functions as a chloride channel and controls the regulation of				
AC	other transport pathways. Mutations in this gene cause the				
PJ	autosomal recessive disorder, cystic fibrosis (CF) and congenital				
VR	bilateral aplasia of the vas deferens (CAVD). Alternative splice				
DP	variants have been described, many of which result from mutations				
	in the CFTR gene.				
	Reviewed				
COMPLETENESS	full length.				
COMMENT	<u>REFSEQ:</u> This reference sequence was derived from M55131.				
PROVISIONAL	RefSeq: This is a provisional reference sequence record that has not yet been subject to human review. The final curated reference sequence record may be somewhat different from this one.				

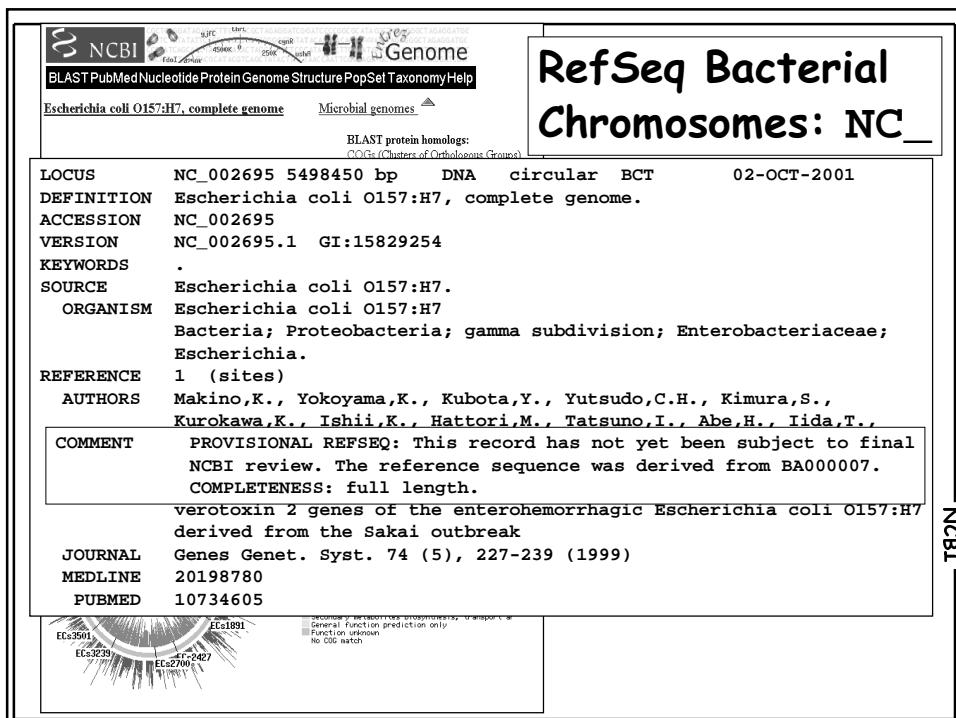
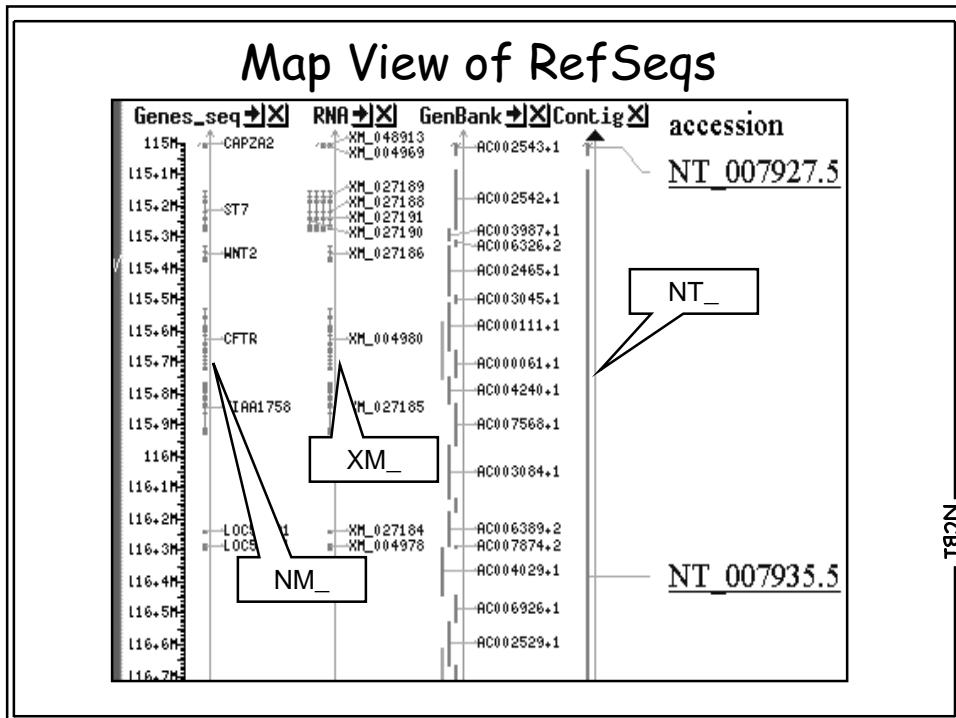
## The Draft Human Genome



N  
CBI

### RefSeq Human Contig: NT\_

LOCUS	NT_007935 1888399 bp DNA	CON	16-NOV-2000
DEFINITION	Homo sapiens chromosome 7 working draft sequence segment, complete sequence.		
ACCESSION	NT_007935		
mRNA	complement(join(1255889..1257642,1258986..1259091,		
CONTIG	join( <u>AC073042.3:1155..2680</u> ,gap(100), <u>AC074390.2:119526..151445</u> , <u>gap(100)</u> , <u>AC074390.2:1..5245</u> ,gap(100), complement( <u>AC074390.2:17705..23645</u> ),gap(100), <u>AC074390.2:97658..119425</u> , <u>AC073042.3:106479..121155</u> , <u>AC074390.2:164226..165036</u> , <u>AC073042.3:70628..79503</u> ,gap(100), <u>AC073042.3:4627..6382</u> ,gap(100), <u>AC073042.3:2781..4526</u> ,gap(100), complement( <u>AC073042.3:183627..29083</u> ),gap(100), <u>AC073042.3:79604..88622</u> ,gap(100), <u>AC073042.3:139234..160437</u> , gap(100),complement( <u>AC073042.3:6483..7909</u> ),gap(100), complement( <u>AC073042.3:39354..45372</u> ),gap(100), complement( <u>AC073042.3:21461..24064</u> )		
	<u>AC074390.2:156347..160294</u> ,gap(100))		<b>Reordering draft sequence</b>
	complement( <u>AC074390.2:5346..10750</u> ),gap(100), complement( <u>AC074390.2:153911..156246</u> ),gap(100), complement( <u>AC074390.2:23746..32402</u> ),gap(100), complement( <u>AC074390.2:151546..153810</u> ),gap(100), complement( <u>AC074390.2:57277..75275</u> ),gap(100), complement( <u>AC074390.2:75376..97557</u> ),gap(100), /gene="CFTR"		



1: NC\_003076  
Arabidopsis thaliana chromosome 5, complete seq  
gi|18426882|ref|NC\_003076.2|[18426882]

2: NC\_003074

LOCUS	NC_003076	26689408 bp	DNA	linear	PLN	10-JAN-2002
DEFINITION	Arabidopsis thaliana chromosome 5, complete sequence.					
ACCESSION	NC_003076					
VERSION	NC_003076.2	GI:18426882				
KEYWORDS	HTG.					
SOURCE	thale cress.					
ORGANISM	Arabidopsis thaliana Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; eudicotyledons; core eudicots; Rosidae; eurosids II; Brassicales; Brassicaceae; Arabidopsis.					
REFERENCE	1	(bases 1 to 26689408)				
AUTHORS	Town,C.D., Haas,B.J., Wu,D., Maiti,R., Hannick,L.I., Chan,A.P., Tallon,L.J., Rooney,T., Utterback,T.R., VanAken,S.E., Feldblyum,T.V., White,O. and Fraser,C.M.					
TITLE	Arabidopsis thaliana chromosome 5 genomic sequence					
JOURNAL	Unpublished					
REFERENCE	2	(bases 1 to 26689408)				
AUTHORS	Town,C.D. and Kaul,S.					
TITLE	Direct Submission					
JOURNAL	Submitted (10-JAN-2002) to the Institute for Genomic Research, 9712 Medical Center Dr, Rockville, MD 20850, USA, cdtown@tigr.org					
COMMENT	PROVISIONAL REFSEQ: This record has not yet been subject to final NCBI review. The reference sequence was derived from AE502093. On Jan 30, 2002 this sequence version replaced gi:15237134. Address all correspondence to:at@tigr.org					

NCBI

**RefSeq Plant  
Chromosomes: NC\_**

**Provisional  
record**

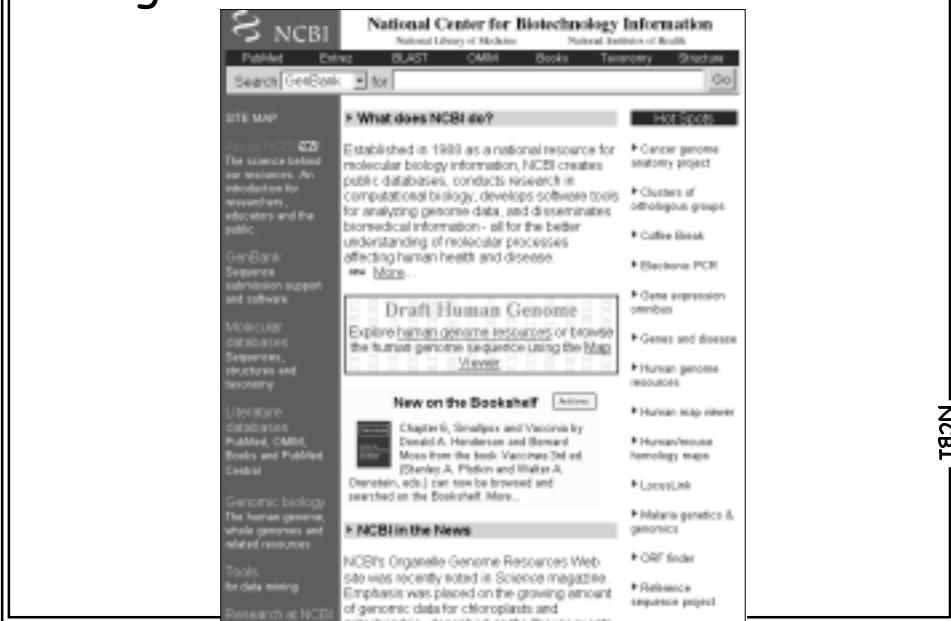
## Other NCBI Derivative Databases

**UniGene** - gene oriented expressed sequence clusters

**LocusLink** - central resource and interface for known genes

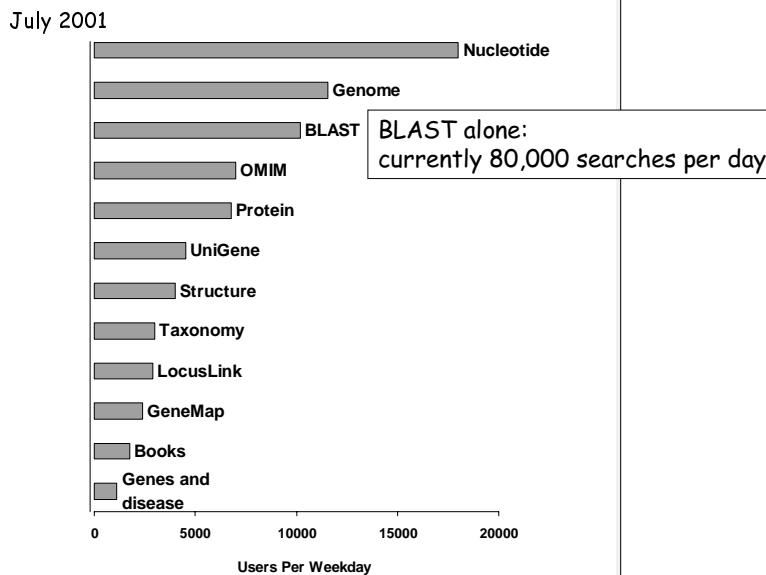
NCBI

## Integrated WWW Access: BLAST and Entrez



NCBI

## Some Web Statistics

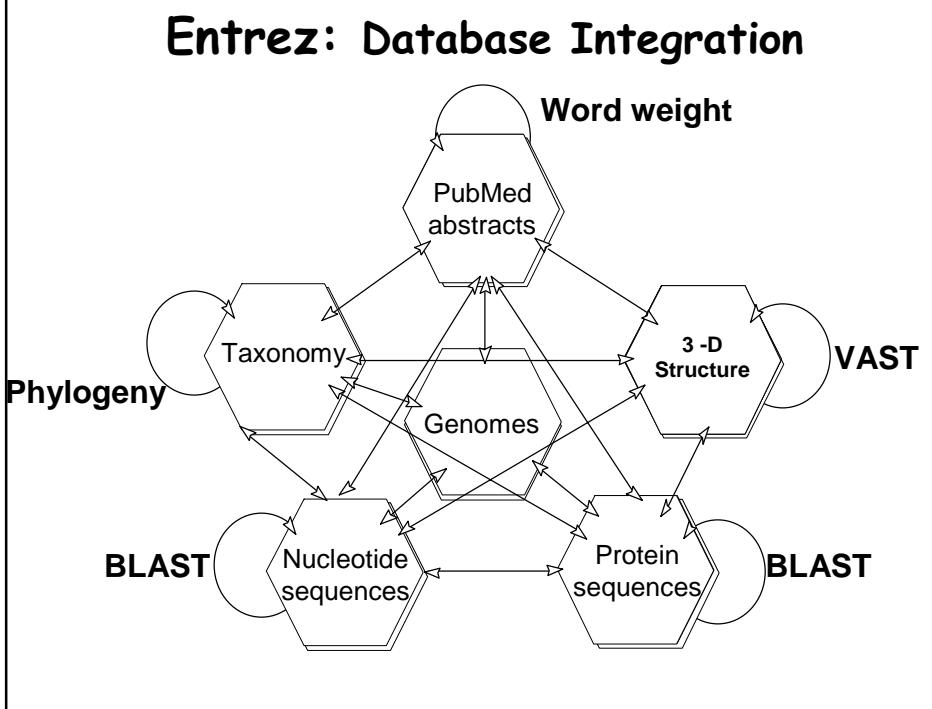


NCBI

## Using Entrez

An integrated database  
search and retrieval system

NCBI



**WWW Entrez**

PubMed: The biomedical literature  
• All of MEDLINE plus others  
• Abstracts  
• Links to online Journals

Nucleotide: GenBank, EMBL, DDBJ  
RefSeq, PDB

Protein sequence database: GenBank, DDBJ, EMBL translations

Structure: three-dimensional protein structures  
PDB, PIR, SWISS-PROT, PRF, RefSeq

NCBI's MMDB - derived from PDB

Genome: ~~comparative genomics~~  
Reference Genomes:  
Graphical views, assembled sequence and mapping data

PopSet: population study data sets

OMIM: Online Mendelian Inheritance in Man

Taxonomy: organisms in GenBank

Books: online books

ProbeSet: Gene Expression Omnibus (GEO)

3D Domains: domains from Entrez Structure

NCBI

**Database Searching with Entrez**

- ◆ Using limits and field restriction to find plant g6pdh
- ◆ Linking and neighboring with g6pdh

NCBI

**Entrez Nucleotides**

The screenshot shows the NCBI Entrez Nucleotides search interface. At the top, there's a decorative graphic of nucleotides (A, T, C, G) and the NCBI logo. Below the logo is a navigation bar with links: PubMed, Nucleotide, Protein, Genome, Structure, PopSet, Taxonomy, OMIM, Books. The 'Nucleotide' link is highlighted. A search bar contains the text 'glucose 6 phosphate dehydrogenase'. Below the search bar are buttons for 'Go' and 'Clear', and links for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. The main area below the search bar is currently empty, indicating no results have been loaded.

**Document Summaries:**

glucose 6 phosphate dehydrogenase[All Fields]

Show: 20 ▾ Items 1-20 of 705 Page 1 of 36 Select page: 1 2 3 4 5 6 7 8 9 10 ×

- 1: NC\_001146** Related Sequences, Protein, PubMed, Taxonomy  
Saccharomyces cerevisiae chromosome XIV, complete chromosome sequence  
gi|6323989|ref|NC\_001146.1|[6323989]
- 2: BM986357** Taxonomy  
EST531265 Rat gene index, normalized rat, norvegicus Rattus norvegicus cDNA clone R:GIAE85 5' end similar to glucose-6-phosphate dehydrogenase, mRNA sequence  
gi|19705746|gb|BM986357.1|[19705746]
- 3: NC\_003424** Protein, PubMed, Taxonomy  
Schizosaccharomyces pombe chromosome I, complete sequence  
gi|19113674|ref|NC\_003424.1|[19113674]
- 4: AJ422253** Related Sequences, Protein, Taxonomy  
Streptococcus pneumoniae partial gdh gene for glucose-6-phosphate dehydrogenase, strain 1338/1996  
gi|19351928|emb|AJ422253.1|SPN422253[19351928]
- 5: NC\_002678** Related Sequences, Protein, PubMed, Taxonomy  
Mesorhizobium loti, complete genome  
gi|13470324|ref|NC\_002678.1|[13470324]
- 6: AL136851** Related Sequences, OMIM, Protein, Taxonomy, UniSTS

**Entrez Nucleotides: Limits & Preview/Index**

The screenshot shows the NCBI Entrez Nucleotides search interface. At the top, there is a search bar with the text "Search Nucleotide for glucose 6 phosphate dehydrogenase". Below the search bar are several buttons: "Limits", "Preview/Index", "History", "Clipboard", and "Details". Two white arrows point upwards from the bottom of the page towards these buttons.

**Entrez**

The screenshot shows the NCBI Entrez search interface. On the left, there is a "Limited to:" dropdown menu with "All Fields" selected. An arrow points to this dropdown from the top of the page. To the right of the dropdown are various search filters: "Accession", "Author Name", "EC/RN Number", "Feature key", "Filter", "Gene Name", "Issue", "Journal Name", "Keyword", "Modification Date", "Organism", "Page Number", "Primary Accession", "Properties", "Protein Name", "Publication Date", "SeqID String", "Sequence Length", "Substance Name", "Text Word", "Title Word", "Uid", and "Volume".

**Entrez Nucleotides: Limits**

NCBI

PubMed Nucleotide Protein Genome Structure PopSet Taxonomy OMIM Books

Search  for **glucose 6 phosphate dehydrogenase**

Go  Clear

**Title == Definition**

**Limited to:**

Title  Exclude Bulk Sequences

exclude ESTs  exclude STSs  include GSS  exclude working draft  
 exclude patents  exclude all of the above

mRNA  Genomic DNA/RNA  Segmented Sequences  
 Only from  Modification Date  Nuclear gene  
 Modification  To  
 Use the mRNA molecule type and day are optional.

NCBI

**Document Summaries: Limits**

Field: Title, Limit: exclude all of the above, Genomic DNA/RNA, mRNA

Display  Summary  Save  Ted  Clip Add

Show  20  Items 1-20 of 61 Page 1 of 4 Select page: 1 2 3 4

1: AY078072 Related Sequences, Protein, Taxonomy  
*Oryza sativa glucose-6-phosphate dehydrogenase (g6pdh) mRNA, complete cds*  
`gi|19071786|gb|AY078072.1|[19071786]`

2: XM\_049337 Related Sequences, Protein, Taxonomy  
*Homo sapiens glucose-6-phosphate dehydrogenase (G6PD), mRNA*  
`gi|14768486|ref|XM_049337.1|[14768486]`

3: NM\_019468 Related Sequences, Protein, PubMed, Taxonomy, UniSTS  
*Mus musculus glucose-6-phosphate dehydrogenase 2 (G6pd2), mRNA*  
`gi|13937388|ref|NM_019468.1|[13937388]`

4: NM\_008052 Related Sequences, Protein, PubMed, Taxonomy, UniSTS  
*Mus musculus glucose-6-phosphate dehydrogenase X-linked (G6pdh), mRNA*  
`gi|6996916|ref|NM_008052.1|[6996916]`

5: U169038 Protein, Taxonomy  
*Hyalophora cecropia glucose-6-phosphate dehydrogenase mRNA, partial cds*  
`gi|6066259|gb|U09406.1|HCU169038|6066259`

6: AY056232 Related Sequences, Protein, Taxonomy  
*Arabidopsis thaliana putative Glucose-6-phosphate dehydrogenase (F14J9.8) mRNA, complete cds*  
`gi|15810386|gb|AY056232.1|[15810386]`

7: BC000337 Related Sequences, OMIM, Protein, Taxonomy

NCBI

## Adding Terms: Preview/Index

[Accession](#)  
[All Fields](#)  
[Author Name](#)

**Add Term(s) to Query or View Index:**

- Enter a term in the text box; use the pull-down menu to specify a search field.
- Click Preview to add terms to the query box and see the number of search results, or click Index to view terms within a field.
- Multiple terms selected from Index will be ORed; click AND to add to search.

Organism

Click **AND** | **OR** | **NOT** to add terms selected from Index to the query box.

green      nts(2774214)  
green      son berry(1)  
green      ccelain crab(11)  
green pygmy goose(1)  
green rice leafhopper(1)  
green ringtail possum(1)  
green rock lobster(2)  
green sea mat(3)  
green seaturtle(105)  
green striped grasshopper(5)

Up      Down

## Plant cytosolic g6pdh mRNAs

Organism Summary Save Text Clip Add

Show 20 21 of 21 One page

<input type="checkbox"/> 1: AJ0078071 Oryza sativa glucose-6-phosphate dehydrogenase mRNA gi 19071796 gb AJ0078071.1  [19071]	<input type="checkbox"/> 15: AJ000183 Spinacia oleracea mRNA for glucose-6-phosphate dehydrogenase, clone O30A5 gi 2276345 emb AJ000183.1  S0000183[2276345]
<input type="checkbox"/> 2: AJ0056232 Arabidopsis thaliana putative glucose-6-phosphate dehydrogenase mRNA gi 1891036 gb AJ0056232.1  [1891036]	<input type="checkbox"/> 16: AJ000182 Spinacia oleracea mRNA for glucose-6-phosphate dehydrogenase, clones O28FA14 & O30A4 gi 2276343 emb AJ000182.1  S0000182[2276343]
<input type="checkbox"/> 3: AB029456 Triticum aestivum g6pdh mRNA, clone Tagp41 gi 8918935 gb AB029456.1  AB029	<input type="checkbox"/> 17: X99405 N.tabacum mRNA for chloroplast glucose-6-phosphate dehydrogenase gi 1480343 emb X99405.1  NTGGPD 1480343
<input type="checkbox"/> 4: AB029457 Triticum aestivum g6pdh mRNA, clone Tagp42 gi 8918936 gb AB029457.1  AB029	<input type="checkbox"/> 18: U18238 Medicago sativa glucose-6-phosphate dehydrogenase mRNA, complete cds gi 603218 gb U18238.1  MSU18238[603218]
<input type="checkbox"/> 5: AB029458 Triticum aestivum g6pdh mRNA, clone Tagp43 gi 8918937 gb AB029458.1  AB029	<input type="checkbox"/> 19: X83923 S.tuberosum mRNA for glucose-6-phosphate dehydrogenase gi 1197384 emb X83923.1  STG6PDHPI 1197384
<input type="checkbox"/> 6: AJ279988 Beta vulgaris partial cDNA for G6PDH gi 672146 gb AJ279988.1  [672146]	<input type="checkbox"/> 20: X74421 S.tuberosum mRNA for glucose-6-phosphate dehydrogenase gi 471344 emb X74421.1  STG6PDH 471344
<input type="checkbox"/> 7: AJ001770 Nicotiana tabacum mRNA for g6pdh gi 3021599 nuc AJ001770.1  NTG6PDHE5[1166404]	<input type="checkbox"/> 21: X84229 A.thaliana mRNA for glucose-6-phosphate dehydrogenase (clone E5) gi 1166404 emb X84229.1  ATG6PDHE5[1166404]
<input type="checkbox"/> 8: AJ001795 Nicotiana tabacum mRNA, cytosolic glucose-6-phosphate dehydrogenase, TIC36	Related Sequences, Protein, PubMed, Taxonomy

**Plant cytosolic *g6pdh* mRNAs**

1: AY078072  
Oryza sativa glucose-6-phosphate  
dehydrogenase mRNA  
gi|19071786|gb|AY078072.1||19071786|

2: AY056232  
Arabidopsis thaliana putative Glucose-6-phosphate dehydrogenase mRNA  
gi|15810386|gb|AY056232.1||15810386|

3: AB029456  
Triticum aestivum g6pdh mRNA  
Elongation Factor 4D  
gi|8911859|gb|AB029456.1||AB029456|

4: AB029457  
Triticum aestivum g6pdh mRNA  
Elongation Factor 4D  
gi|8911860|gb|AB029457.1||AB029457|

5: AB029451  
Triticum aestivum g6pdh mRNA  
Elongation Factor 4D  
gi|8911851|gb|AB029454.1||AB029451|

6: AJ279688  
Beta-tubulin partial mRNA for C.  
gi|723464|emb|AJ279688.1||AJ279688|

7: AJ001770  
Nicotiana tabacum mRNA for g6pdh  
gi|0421599|ncbi|AJ001770.1||NTU001770|

8: AJ001769  
Nicotiana tabacum mRNA cytosolic glucose-6-phosphate dehydrogenase TC06  
gi|0421598|ncbi|AJ001769.1||NTU001769|

**Summary**  
**Brief**  
**GenBank**  
**ASN.1**  
**FASTA**  
**GI list**  
**LinkOut**  
**PubMed Links**  
**Protein Links**  
**Nucleotide Neighbors**  
**PopSet Links**  
**Structure Links**  
**Genome Links**  
**Taxonomy Links**  
**OMIM Links**

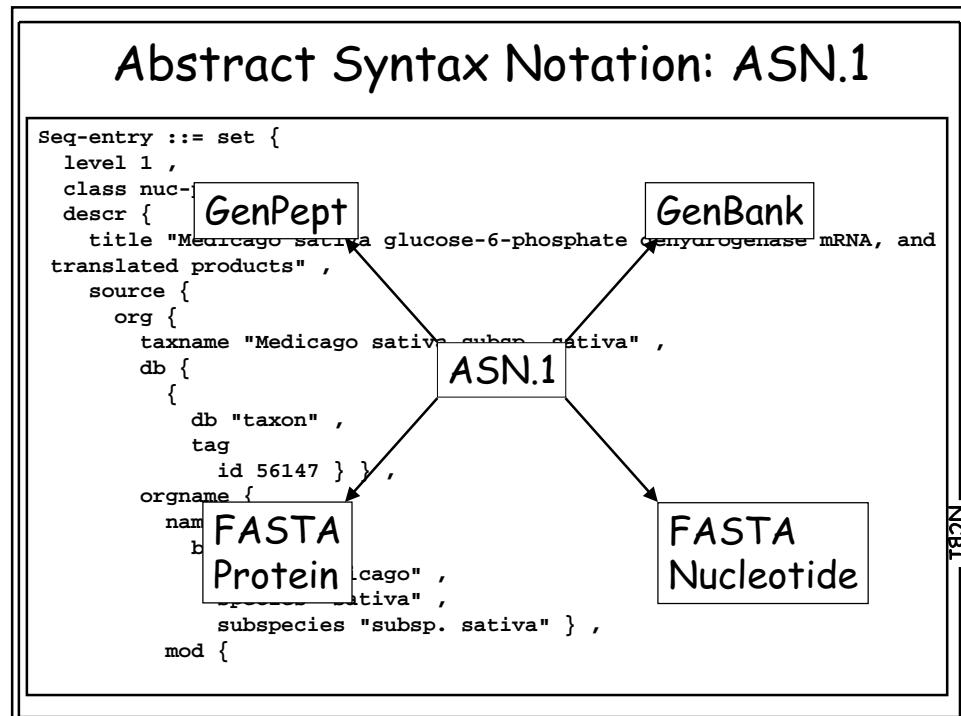
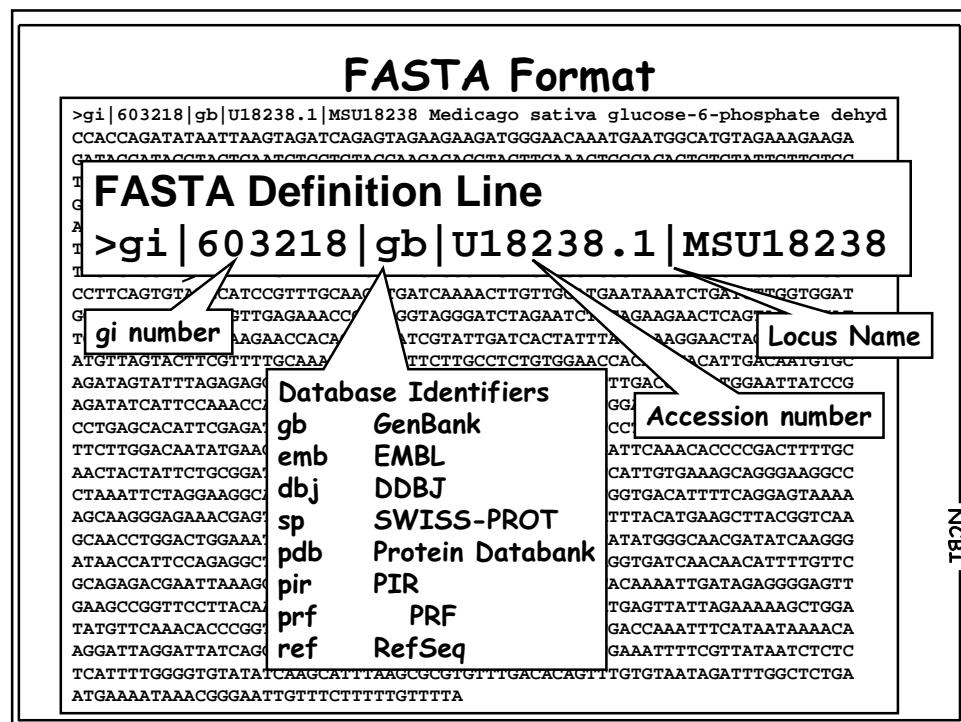
**Entrez GenBank / GenPept**

1: U18238: Medicago sativa g...[gi:603218]

**LOCUS** MSU18238  
**DEFINITION** Medicago sativa glucose-6-phosphate dehydrogenase gene  
**ACCESSION** U18238  
**VERSION** U18238.1 GI:603218  
**KEYWORDS** .  
**SOURCE** alfalfa.  
**ORGANISM** Medicago sativa subsp. sativa Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophytina; Magnoliopsida; Rosidae; eurosids I; Fabales; Fabaceae; Papilionoideae; Trifolieae; Medicago.  
**REFERENCE** 1 (bases 1 to 1785)  
**AUTHORS** Fahrendorf,T., Ni,W.,  
**TITLE** Stress responses in alfalfa at the onset of transcriptional activation  
**JOURNAL** Plant Mol. Biol. 28 (1995) 764-769  
**MEDLINE** 95367649  
**PUBMED** 7640360  
**REFERENCE** 2 (bases 1 to 1785)  
**AUTHORS** Fahrendorf,T.  
**TITLE** Direct Submission  
**JOURNAL** Submitted (06-DEC-1995)  
**FOUNDATION** Foundation, Plant Bioinformatics  
**LOCATION** Ardmore, OK 73402, US  
**FEATURES** Location/Qualifiers  
**source** 1..1785  
**/organism="Medicago sativa"**  
**/cultivar="cv"**  
**/sub\_species="cv"**  
**/db\_xref="taxon:33430"**  
**/clone="G10"**  
**/tissue\_type="cell suspension culture"**  
**CDS**  
**/db\_xref="GenBank:U18238.1:38..1585"**  
**/translation="KTFPAFLFHLYQKELLPPDEEVHIFGYAKRKRISDEELRNRLKSILVPERGASPRQLDDVS"**

1: AAB41552: glucose-6-phosphat...[gi:603219]

**LOCUS** AAB41552 515 aa linear PLN 30-JAN-1997  
**DEFINITION** glucose-6-phosphate dehydrogenase.  
**ACCESSION** AAB41552  
**PID** g603219  
**VERSION** AAB41552.1 GI:603219  
**DBSOURCE** locus MSU18238 accession U18238.1  
**KEYWORDS** .  
**SOURCE** Medicago sativa subsp. sativa  
**ORGANISM** Medicago sativa subsp. sativa Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophytina; Magnoliopsida; Rosidae; eurosids I; Fabales; Fabaceae; Papilionoideae; Trifolieae; Medicago.  
**REFERENCE** 1 (residues 1 to 515)  
**AUTHORS** Fahrendorf,T., Ni,W., Shorosh,B.S. and Dixon,R.A.  
**TITLE** Stress responses in alfalfa (Medicago sativa L.) XIX. Tissue specific expression of the glucose-6-phosphate dehydrogenase gene  
**PROTEIN** 1..515  
**/product="glucose-6-phosphate dehydrogenase"**  
**CDS** 1..515  
**/coded\_by="U18238.1:38..1585"**  
**ORIGIN** 1 mgtnewhver rdsigtespv arekvletgtl sivvlgasgd lakkktfpal fhlykqellp 61 pdevhifgya rskisddel nklrsylvpe kgaspkqidd vskflqlvky vsgpydsedq 121 frlldkeise heylksnkg srrlfyfl psppvpsvc miktccmns diggvtvv 181 ekpfgrdles aeelstqige lfeepqiyai dhylgkelvq nmlylrfanr fflpliwhnh 241 idnvqivfrie dftgdgrgyy fdqygiirdi ipnllqvlc liamekpvs1 kpehirdek 301 kyvlesvlpir ddevvlqgye gytdpptvpd dsnptffatt ilrihnerwe gvpfivkagk 361 alnsrkaeir vqfkdvpgdi frskkkqrne fvirlgpsea iymklvkgp glemsavqse 421 ldlsyqrgya gitipeayer lildtirgqd qhvrrdelk aswqiftpl1 hkidrgelkp 481 vpynpgsrgp aeadelleka gyvtpgyiwi pptl



# NCBI Toolbox

**Toolbox Sources**

```

/*
 *      asn2ff.c
 *      convert an ASN.1 entry to flat file format. using the FFPrintArray.
 *
***** */
#include <accent.h>
#include "asn2ff.h"
#include "asn2ff.h"
#include "ffprint.h"
#include <subutai.h>
#include <objall.h>
#include <objcod.h>
#include <lsqfet.h>
#include <explor.h>

#ifndef ENABLE_IDL
#include <accidl.h>
#endif

FILE *fp1;

```

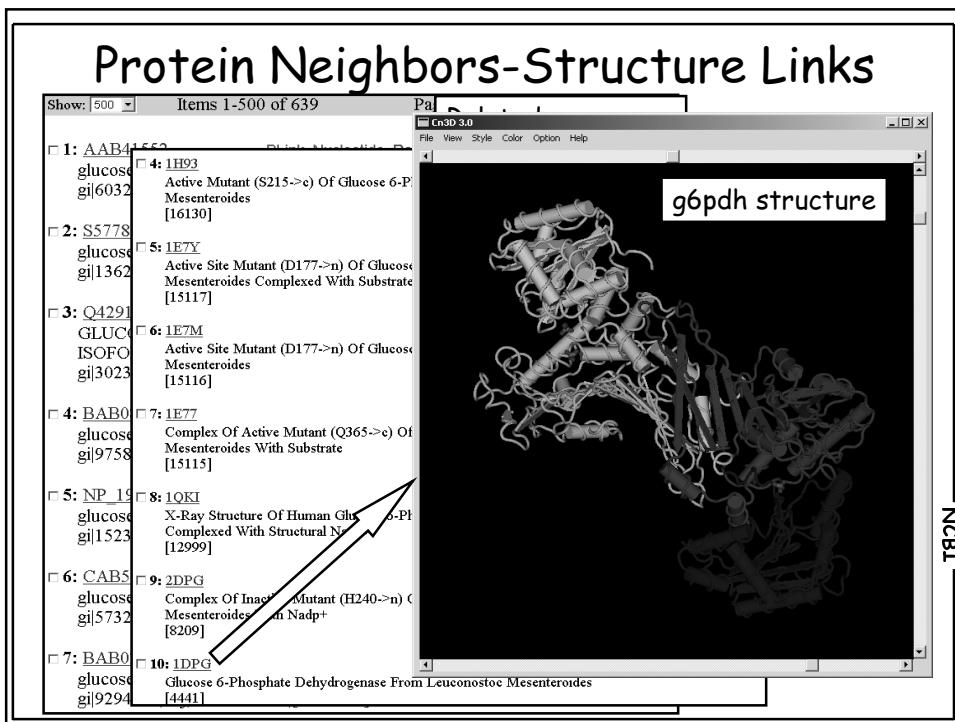
**ftp> open ftp.ncbi.nih.gov**

**ftp> cd toolbox**

**ftp> cd ncbi\_tools**

Args myarg  
**ftp://ftp.ncbi.nlm.gov/toolbox/ncbi\_tools**  
 {"Input asnfile in binary mode","F",NULL,NULL,TRUE,'b',ARG\_BOOLEAN,0,0,0,NULL},  
 {"Output Filename","stdout", NULL,NULL,TRUE,'o',ARG\_FILE\_OUT,0,0,0,NULL},  
 {"Show Sequence?","T", NULL ,NULL ,TRUE,'h',ARG\_BOOLEAN,0,0,0,NULL},

NCBI



NCBI

## Advanced Neighbors: BLink

1: [AAB41552](#) BLink, Nucleotide, Related Sequences, PubMed, Taxonomy  
 glucose-6-phosphate dehydrogenase [Medicago sativa subsp. sativa]  
 gi|603219|gb|AAB41552.1|[60 19]

NCBI

## BLink

BLAST    PubMed    Nucleotide    Protein    Genome    Structure    Taxonomy    Help

Query: gi|603219 glucose-6-phosphate dehydrogenase

Matching db: 1362052 3027915

515 aa

SCORE	P	ACCESSION	GI	N	ORGANISM
2250	14	gi 603219	603219	1	Ambidioecia thalassina
		350	3	AAL18424	16416802 - 1 Aedeomyia squamipennis
		347	3	AAL18427	16416808 - 1 Anopheles stephensi
		340	3	AAL18435	16416824 - 1 Anopheles coustani
		339	3	AAL18426	16416806 - 1 Toxorhynchites amboinensis
		337	3	AAL18443	16416840 - 1 Anopheles albimanus
		336	3	AAL18425	16416804 - 1 Uranotaenia sapphirina
		335	3	AAL18442	16416838 - 1 Anopheles albitarsis
		332	3	AAB02777	606624 - 1 Ctenolepisma lineata
		331	3	AAL18439	16416832 - 1 Bironella gracilis
		330	3	AAL18432	16416818 - 1 Anopheles mattogrossensis
		327	3	AAL18441	16416836 - 1 Aedimyces subalbatus
					Anopheles quadrimaculatus
					Anopheles bellator
					Anopheles cruzii
		278	3	AAB02786	606642 - 1 Salina tristani
		318	3	AAB02787	606644 - 1 Zerene cesonia
		317	3	AAL18431	16416816 - 1 Anopheles neivai
		316	3	AAL18430	16416814 - 1 Orthopodomyia alba
		313	3	AAB02781	606632 - 1 Sarcophaga bullata
		304	3	AAB02778	606626 - 1 Culex pipiens
		295	3	BAA78540	4996576 - 1 Bactrocera dorsalis
		293	3	BAA78545	4996586 - 1 Bactrocera scutellata
		287	3	BAA78548	4996592 - 1 Bactrocera kandiensis
		286	3	CAB61333	6453559 - 1 Laminaria digitata
		286	3	BAA78543	4996582 - 1 Bactrocera cucurbitae

NCBI

# PubMed Link

[1: Plant Mol Biol 1995 Aug;28\(5\):885-900](#)

Related Articles, Nucleotide, Protein, [New Books](#), LinkOut

**Stress responses in alfalfa (*Medicago sativa L.*) XIX. Transcriptional activation of oxidative pentose phosphate pathway genes at the onset of the isoflavanoid phytoalexin response.**

Fahrendorf T, Ni W, Shorrosh BS, Dixon RA.

Plant Biology Division, Samuel Roberts Noble Foundation, Ardmore, OK 73401, USA.

We have isolated cDNA clones encoding the pentose phosphate pathway enzymes 6-phosphogluconate dehydrogenase (6PGD, EC 1.1.1.44) and glucose 6-phosphate dehydrogenase (G6PDH, EC 1.1.1.49) from alfalfa (*Medicago sativa L.*). These exhibit extensive nucleotide and amino acid sequence similarity to the corresponding genes from bacteria, *Drosophila* and mammals. Transcripts encoding both enzymes are expressed at high levels in roots and nodules. Exposure of alfalfa suspension cells to an elicitor from yeast cell walls results in co-ordinated increases in transcription rates for both genes, followed by increased steady state transcript levels but only slightly increased extractable enzyme activities, at the onset of accumulation of isoflavanoid phytoalexins. Levels of NADPH and NADP remain relatively constant in alfalfa cells following elicitation. The rapid transcriptional activation of 6PGDH and G6PDH does not therefore appear to be a response to altered pyridine nucleotide redox state. These genes appear to respond to early events in elicitor-mediated signalling rather than to subsequent elicitor-induced changes in secondary metabolism. Hydrogen peroxide, a potential signal for elicitation of anti-oxidative genes in biologically stressed plant cells did not induce 6PGDH or G6PDH transcripts or enzymatic activity.

## Online Books

**MOLECULAR BIOLOGY OF THE CELL**

**NCBI**

**Navigation**

- 7. Recombinant DNA Technology
- Introduction
- The Fragmentation, Separation, and Sequencing of DNA Molecules
- Nucleic Acid Hybridization
- DNA Cloning
- DNA Engineering
- References

**Search**

This book    All books

**The Tax Browser**

The NCBI taxonomy browser interface. The search bar shows 'Search for [Taxonomy] for Fragaria'. The results list includes:

- Fragaria chiloensis [14]
  - Fragaria chiloensis subsp. lucida [2]
  - Fragaria chiloensis subsp. pacifica [2]
- Fragaria daltoniana [2]
- Fragaria gracilis [2]
- Fragaria iinumae [92]
- Fragaria moschata [21]
- Fragaria nilgerrensis [21]
- Fragaria nippponica [2]
- Fragaria rubicunda [86]
- Fragaria orientalis [6]
- Fragaria pentaphylla [4]
- Fragaria vesca [5925]
  - Fragaria vesca subsp. americana [2]
  - Fragaria vesca subsp. bracteata [22]
    - Fragaria vesca subsp. vesca [106]
      - Fragaria vesca f. alba [2]
  - Fragaria virginiana [14]
    - Fragaria virginiana subsp. glauca [2]
    - Fragaria virginiana subsp. platypetala [4]
    - Fragaria virginiana subsp. virginiana [6]
  - Fragaria viridis [51]
  - Fragaria x ananassa (strawberry) [160]
  - Fragaria sp. CFRA 538 [2]

**TaxBrowser: Rose Family**

Rosaceae

**Taxonomy ID: 3745**  
Rank: family  
Genetic code: Translation table  
Mitochondrial genetic code: T

Other names:  
**rose family** [common name]  
Lineage (abbreviated)  
Eukaryota; Viridiplantae; Magnoliophyta; eudicots

Entrez Nucleotide Help | FAQ

Batch Entrez: Upload a file of GI or accession numbers to retrieve sequences

Check sequence revision history  
How to create www links to Entrez  
LinkOut  
Cubby  
Related

Search for Genes LocusLink provides curated information for human, fruit fly, mouse, rat, and zebrafish

Nucleotide (6509) Protein

NCBI Taxonomy Browser

Search for [Nucleotide] for txid3745[Organism]

Display: Summary Save Text Clip Add

Show: 20 Items 1-20 of 6509 Page 1 of 326 Select page: 1 2 3 4 5 6 7 8 9 10 >

1: AF454000 Protein, Taxonomy  
Prunus dulcis cultivar Anxaneta RNase (S) gene, S2 allele, partial cds  
gi|18092541|gb|AF454000.1|[18092541]

2: AY061960 Related Sequences, Protein, Taxonomy  
Prunus dulcis cultivar Ferragnes RNase (S) gene, S1 allele, partial cds  
gi|17266291|gb|AY061960.1|[17266291]

3: AF323507 Protein, Taxonomy  
Malus x domestica sorbitol dehydrogenase (SDH4) mRNA, complete cds  
gi|17225199|gb|AF323507.1|[17225199]

4: AF323506 Protein, Taxonomy  
Malus x domestica sorbitol dehydrogenase (SDH3) mRNA, complete cds  
gi|17225197|gb|AF323506.1|[17225197]

5: AF323505 Protein, Taxonomy  
Malus x domestica sorbitol dehydrogenase (SDH2) mRNA, complete cds  
gi|17225195|gb|AF323505.1|[17225195]

6: AF323504 Related Sequences, Protein, Taxonomy  
Malus x domestica sorbitol dehydrogenase (SDH1) mRNA, complete cds

## Entrez Structures

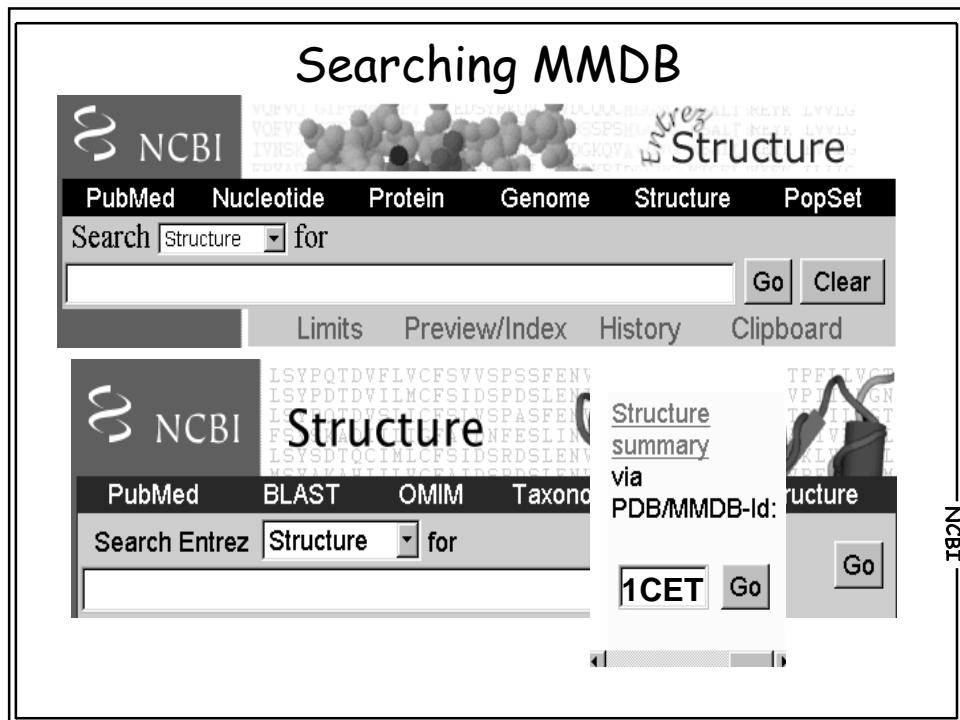
Molecular Modeling Database (MMDB)  
and Ch3D

NCBI

### MMDB: Molecular Modeling Data Base

- Derived from experimentally determined PDB records
- Value added to PDB records including:
  - Addition of explicit chemical graph information
  - Validation
  - Inclusion of Taxonomy, Citation, and other information
  - Conversion to ASN.1 data description language
- Structure neighbors determined by  
Vector Alignment Search Tool (VAST)

NCBI



## Structure Summary

MMDB Id: 9873 PDB Id: 1CET

**Protein Chains:** A

**MEDLINE:** [PubMed](#)

**Taxonomy:** A [Plasmodium falciparum](#)

**PDB Authors:** J.A.Read, K.W.Wilkinson, R.T.

**BLAST neighbors**

**Sequence Neighbors:** A

**Structure Neighbors:** A

**VAST neighbors**

[View / Save Structure](#)

**Options:**

- Launch Viewer
- See File
- Save File

**Viewer:**

- Cn3D (asn.1)
- Cn3D v1.0 (asn.1)
- Mage
- RasMol (PDB)

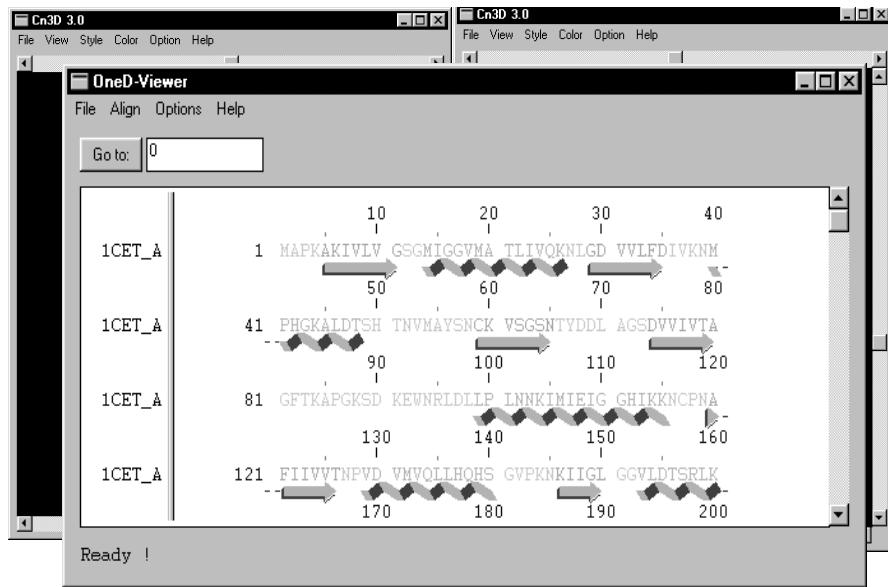
**Complexity:**

- Cn3D Subset
- Virtual Bond Model
- All Atom Model
- Up to 5 Models
- Up to 10 Models
- All Models

**Cn3D viewer**

[Get Cn3D 3.0!](#)

## Cn3D : Displaying Structures

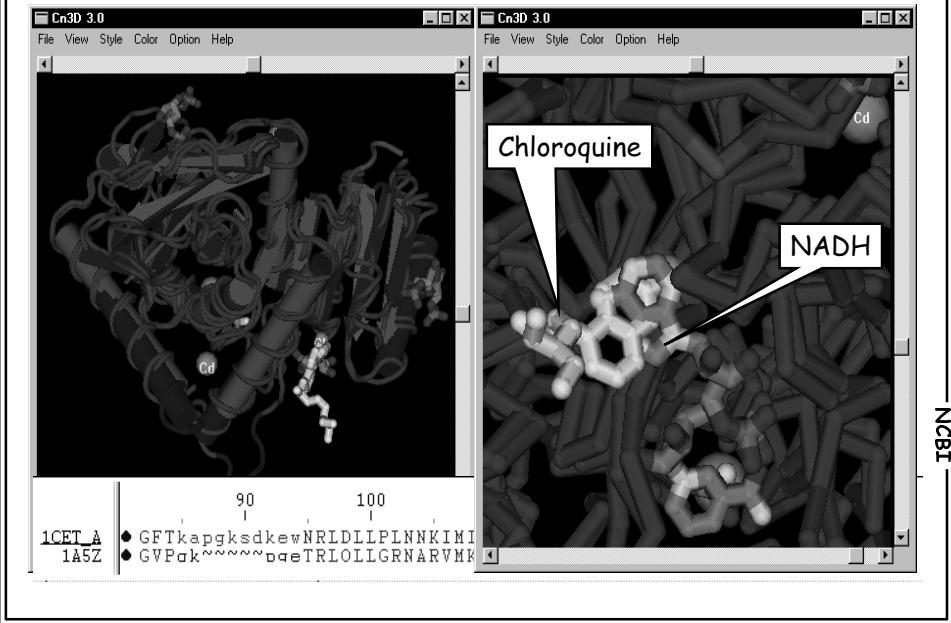


## Structure Neighbors

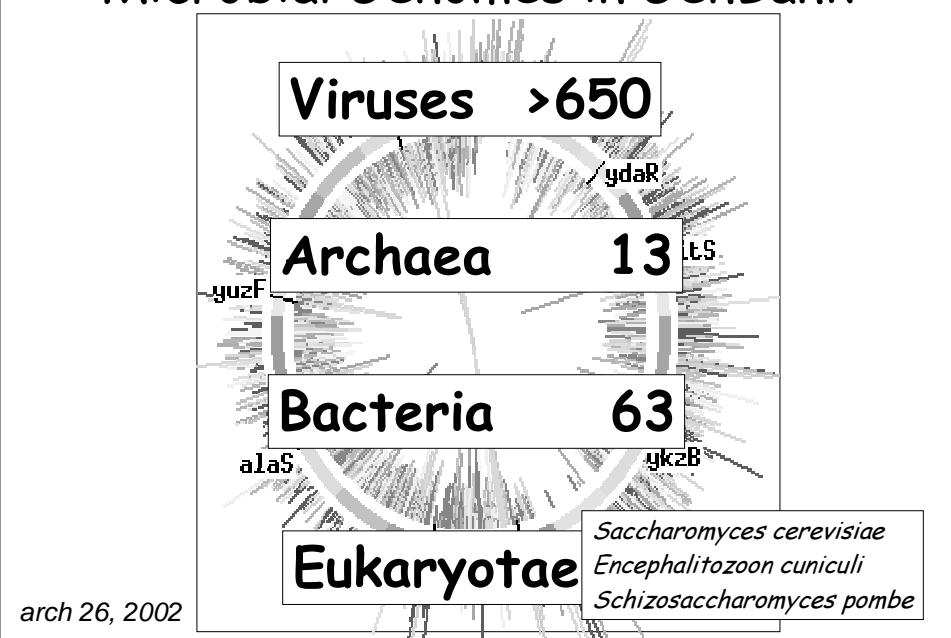
Structure neighbors 1-20 out of 208 displayed. Page 1 of 11.

	PDB	C	D	RMSD	NRES	%Id	Description
<input type="checkbox"/>	<a href="#">1LLD_A</a>			1.8	294	28.9	L-Lactate Dehydrogenase (E.C.1.1.1.27) (T-State) Mutant With Cys 199 Replaced By Ser (C199s) Complex With Nadh
<input type="checkbox"/>	<a href="#">1LDG</a>			0.6	300	100.0	Plasmodium Falciparum L-Lactate Dehydrogenase Complexed With Nadh And Oxamate
<input type="checkbox"/>	<a href="#">1CEQ_A</a>			0.2	304	99.7	Chloroquine Binds In The Cofactor Binding Site Of Plasmodium Falciparum Lactate Dehydrogenase.
<input type="checkbox"/>	<a href="#">6LDH</a>			1.7	298	27.8	M4 Apo-Lactate Dehydrogenase (E.C.1.1.1.27)
<input checked="" type="checkbox"/>	<a href="#">1A5Z</a>			1.6	300	36.3	Lactate Dehydrogenase From Thermotoga Maritima (Tmldh)
<input type="checkbox"/>	<a href="#">1LDN_A</a>			1.6	300	28.7	L-Lactate Dehydrogenase (E.C.1.1.1.27) Complexed With Nadh, Oxamate, And Fructose-1,6-Bisphosphate

## Structural Alignments



## Microbial Genomes in GenBank



**Bacterial Genomes**

Bacteria Complete Genomes Taxonomy / List 50	
Archaeabacteria bacteria	NC_003682 28415
Aquifex aeolicus	NC_003683 2034
Bacillus halodurans	NC_000918 1551
Bacillus subtilis	NC_002570 4282
Bacillus thuringiensis	NC_000964 4214
Borrelia burgdorferi	NC_001318 917
Burkholderia sp. APS	NC_002528 649
Candidatus etus	NC_002163 1643
Candidatus crescens	NC_002696 4016
Chlorobium variabile	NC_002620 1069
Chlorobium variabile	NC_001117 1042
Chlorobium pneumoniae AB39	NC_002179 1228
Chlorobium pneumoniae CBT109	NC_000922 1230
Chlorobium pneumoniae J138	NC_002491 1226
Chlorobium pneumoniae	NC_003030 3848
Dinoroseococcus radiotolerans	NC_001263 2649
	NC_001264 412
Escherichia coli K-12	NC_000913 4639
Escherichia coli O157:H7	NC_002695 5498
Escherichia coli O157:H7 EDL933	NC_002655 5528
Haemophilus influenzae R-4	NC_000007 1830
Helicobacter pylori 26695	NC_000915 1687
Helicobacter pylori 399	NC_000921 1643
Lactococcus lactis subsp. lactis	NC_002662 2385
Mesorhizobium loti	NC_002678 7036
	NC_002478 7036074 bp Sep 10 2001
	NC_002677 3268263 bp Feb 20 2001
	NC_00255 4409394 bp Apr 25 2001
	NC_00062 4411529 bp Sep 7 2001
	NC_00246 588074 bp Jun 8 2001
	NC_0112 814394 bp Apr 2 2001
	NC_0771 963879 bp May 14 1991
	NC_0183 2272351 bp Feb 25 2000
	NC_00203 2184406 bp Mar 30 2000
	NC_002652 2257487 bp Sep 10 2001
	NC_002516 6364403 bp Sep 10 2001
	NC_003103 1268755 bp Apr 26 2001
	NC_00063 1111523 bp Sep 10 2001
	NC_003647 3454135 bp Jul 30 2001
	NC_002158 2878048 bp Jan 1 2001
	NC_002482 2813643 bp Apr 21 2001
	NC_003996 2038615 bp Jul 27 2001
	NC_003628 2164837 bp Jul 25 2001
	NC_002337 3852443 bp Apr 10 2001
	NC_00011 3573479 bp Mar 21 1997
	NC_00053 3868725 bp Sep 10 2001
	NC_00019 1138011 bp Sep 7 2001
	NC_00262 751719 bp Jan 18 2001
	NC_002505 2961149 bp Sep 10 2001
	NC_002506 3072315 bp Sep 10 2001

***M. tuberculosis* Complete Genome**

Protein coding genes distribution map  
To see map locations of genes, click on a region in

[Mycobacterium tuberculosis H37Rv complete genome](#) [Microbial genomes](#)

**BLAST protein homologs:**  
COGs (Clusters of Orthologous Groups)  
3D Structure (Sequences with known structure)  
TaxMap (Sequences grouped by superkingdom)  
TaxPlot (3-way genome comparison)  
CDD(Conserved Domain Database)

**Feature table:**  
Protein coding genes: 3918  
Structural RNAs: 54

**Contributor:** [Sanger Centre](#) See genome at [Sanger Centre](#)  
Download chromosome sequence data from [NCBI FTP site](#)

[NEW BLAST your query sequence against the genome](#) ▶ [NEW BLAST against protein sequences](#) ▶

Organism: [Mycobacterium tuberculosis](#)  
Genetic Code: 11  
Lineage: Eubacteria; Firmicutes; Actinomycetes; Mycobacteria; Mycobacteriaceae; Mycobacterium

## Coding Regions

### Mycobacterium tuberculosis H37Rv complete genome

[Download \[Table\]](#) [\[FASTA protein\]](#) [\[FASTA nucleotide\]](#) from NCBI FTP site

◆ - GenBank record including protein ◇ - DNA region in flatfile format ◆ - DNA and protein in FASTA format

	Location	Strand	Length	PID	Gene	COG	Synonym	Product
◆ ◇ ◆	1..1524	+	508	<a href="#">2808711</a>	<a href="#">dnak</a>	<a href="#">COG0593</a>	Rv0001	<a href="#">dnak</a>
◆ ◇ ◆	2052..3260	+	403	<a href="#">3261513</a>	<a href="#">dnan</a>	<a href="#">COG0592</a>	Rv0002	<a href="#">dnan</a>
◆ ◇ ◆	3280..4437	+	386	<a href="#">1552556</a>	<a href="#">recF</a>	<a href="#">COG1195</a>	Rv0003	<a href="#">recF</a>
◆ ◇ ◆	4434..4997	+	168	<a href="#">1552557</a>	Rv0004			hypothetical protein Rv0004
◆ ◇ ◆	5123..7267	+	715	<a href="#">1552558</a>	<a href="#">gyrB</a>	<a href="#">COG0187</a>	Rv0005	<a href="#">gyrB</a>
◆ ◇ ◆	7302..9818	+	839	<a href="#">1552559</a>	<a href="#">gyrA</a>	<a href="#">COG0188</a>	Rv0006	<a href="#">gyrA</a>
◆ ◇ ◆	9914..10828	+	305	<a href="#">1552560</a>	Rv0007			hypothetical protein Rv0007
◆ ◇ ◆	11874..12311	-	146	<a href="#">1552562</a>	Rv0008c			hypothetical protein Rv0008c
◆ ◇ ◆	12468..13016	+	183	<a href="#">1552563</a>	<a href="#">ppiA</a>	<a href="#">COG0652</a>	Rv0009	<a href="#">ppiA</a>
◆ ◇ ◆	13133..13558	-	142	<a href="#">1552565</a>	Rv0010c			hypothetical protein Rv0010c
◆ ◇ ◆	13714..13995	-	94	<a href="#">1552566</a>	Rv0011c			hypothetical protein Rv0011c
◆ ◇ ◆	14089..14877	+	263	<a href="#">3261646</a>	Rv0012			hypothetical protein Rv0012
◆ ◇ ◆	14914..15612	+	233	<a href="#">1552568</a>	<a href="#">pabA</a>	<a href="#">COG0512</a>	Rv0013	<a href="#">pabA</a>
◆ ◇ ◆	15590..17470	-	627	<a href="#">1552569</a>	<a href="#">pknB</a>	<a href="#">COG0515</a>	Rv0014c	<a href="#">pknB</a>
◆ ◇ ◆	17467..18762	-	432	<a href="#">1552570</a>	<a href="#">pknA</a>	<a href="#">COG0515</a>	Rv0015c	<a href="#">pknA</a>
◆ ◇ ◆	18759..20234	-	492	<a href="#">1552571</a>	<a href="#">pbpA</a>	<a href="#">COG0768</a>	Rv0016c	<a href="#">pbpA</a>
◆ ◇ ◆	20231..21640	-	470	<a href="#">1552572</a>	<a href="#">rodA</a>	<a href="#">COG0772</a>	Rv0017c	<a href="#">rodA</a>
◆ ◇ ◆	21637..23181	-	515	<a href="#">1552573</a>	<a href="#">ppp</a>	<a href="#">COG0631</a>	Rv0018c	<a href="#">ppp</a>
◆ ◇ ◆	23270..23737	-	156	<a href="#">1552574</a>	Rv0019c	<a href="#">COG1716</a>		hypothetical protein Rv0019c
◆ ◇ ◆	23861..25444	-	528	<a href="#">1552575</a>	Rv0020c	<a href="#">COG1716</a>		hypothetical protein Rv0020c
◆ ◇ ◆	25913..26881	-	323	<a href="#">1552576</a>	Rv0021c	<a href="#">COG2070</a>		hypothetical protein Rv0021c
◆ ◇ ◆	27023..27442	-	140	<a href="#">1552577</a>	Rv0022c			hypothetical protein Rv0022c
◆ ◇ ◆	27595..28365	+	257	<a href="#">1552578</a>	Rv0023	<a href="#">COG1396</a>		hypothetical protein Rv0023
◆ ◇ ◆	28323..29027	-	223	<a href="#">1552579</a>	Rv0024	<a href="#">COG0201</a>		hypothetical protein Rv0024

NCBI

## Genome Annotations

### Mycobacterium tuberculosis H37Rv complete genome 818239..868238

Rv0726c fucR Rv0728c xy1B Rv0730 secY adk map\* sigL

31 Archaea 94 Bacteria 68 Metazoa 0 Fungi 6 Plants 0 Viruses 1 Other Eukaryotae

Keep only No filter Cut-Off 100 Select Reset

390 aa

	SCORE	P	ACCESSION	GI	N	ORGANISM
	1146	3	<a href="#">AAK23331</a>	13422699	-	<a href="#">7 Caulobacter crescentus</a>
	1017	2	<a href="#">AAF12736</a>	6465955	-	<a href="#">21 Homo sapiens</a>
	1001	2	<a href="#">BAB26255</a>	12844144	-	<a href="#">13 Mus musculus</a>
	989	3	<a href="#">AAG04135</a>	9946632	-	<a href="#">13 Pseudomonas aeruginosa</a>
	676	7	<a href="#">CAB61663</a>	6468707	-	<a href="#">12 Streptomyces coelicolor A3(2)</a>
	657	4	<a href="#">BAB07517</a>	10176423	-	<a href="#">4 Bacillus halodurans</a>
	653	3	<a href="#">AAK18172</a>	13310129	-	<a href="#">1 Pseudomonas putida</a>
	647	4	<a href="#">CAA89868</a>	853760	-	<a href="#">7 Bacillus subtilis</a>
	644	3	<a href="#">AAF10828</a>	6459002	-	<a href="#">7 Deinococcus radiodurans</a>
	640	7	<a href="#">CAB61531</a>	6465864	-	<a href="#">1 Streptomyces lividans</a>
	597	4	<a href="#">AAA95968</a>	1055219	-	<a href="#">1 Clostridium acetobutylicum</a>

NCBI

## M. tuberculosis vs. E.coli COGS

Code	COGs	Description	Code	COGs	Description
■ J	140	Translation, ribosomal structure and biogenesis	■ J	167	Translation, ribosomal structure and biogenesis
■ K	165	Transcription	■ K	238	Transcription
■ L	166	DNA replication, recombination and repair	■ L	207	DNA replication, recombination and repair
□ D	28	Cell division and chromosome partitioning	□ D	26	Cell division and chromosome partitioning
■ O	86	Posttranslational modification, protein turnover, chaperones	■ O	117	Posttranslational modification, protein turnover, chaperones
□ M	115	Cell envelope biogenesis, outer membrane	□ M	190	Cell envelope biogenesis, outer membrane
□ N	32	Cell motility and secretion	□ N	119	Cell motility and secretion
□ P	115	Inorganic ion transport and metabolism	□ P	184	Inorganic ion transport and metabolism
□ T	92	Signal transduction mechanisms	□ T	140	Signal transduction mechanisms
■ C	192	Energy production and conversion	■ C	266	Energy production and conversion
■ G	117	Carbohydrate transport and metabolism	■ G	324	Carbohydrate transport and metabolism
□ E	200	Amino acid transport and metabolism	□ E	343	Amino acid transport and metabolism
■ F	68	Nucleotide transport and metabolism	■ F	89	Nucleotide transport and metabolism
■ H	111	Coenzyme metabolism	■ H	115	Coenzyme metabolism
■ I	206	Lipid metabolism	■ I	86	Lipid metabolism
□ R	440	General function prediction only	□ K	336	General function prediction only
□ S	188	Function unknown	□ S	290	Function unknown
☒ -	1456	not in COGs	☒ -	1031	not in COGs

NCBI

## Entrez Genomes

NCBI

PubMed Nucleotide Protein Genome Structure PopSet

Search Genome for  Go Clear

Limits Index History Clipboard

Eukaryote Genomes Taxonomy / List

Complete genome

- [ 5] [Anopheles gambiae](#) NEW chromosomes: I, 2, 3,
- [ 5] [Arabidopsis thaliana](#) chromosomes: I, II, III, IV, V,
- [ 6] [Caenorhabditis elegans](#) chromosomes: I, II, III, IV, V, X,
- [ 5] [Drosophila melanogaster](#) chromosomes: 1, 2, 3, 4, Y,
- [ 11] [Encephalitozoon cuniculi genome](#) chromosomes: I, II, III, IV, V, VI, VII, VIII, IX, X, XI,
- [ 3] [Guillardia theta nucleomorph genome](#) chromosomes: 1, 2, 3,
- [ 16] [Saccharomyces cerevisiae](#) chromosomes: I, II, III, IV, V, VI, VII, VIII, IX, X, XI, XII, XIII, XIV, XV, XVI,
- [ 3] [Schizosaccharomyces pombe](#) chromosomes: I, II, III,

NCBI

**The Arabidopsis Map Viewer**

NCBI

PubMed Nucleotide Protein Genome Structure PopSet Taxonomy OMIM Help

Search for At5g44790 on chromosome(s) Find

Show linked entries Help FTP

Entrez Genomes

Prominent organisms

FTP SITE

Related Databases: TAIR TIGR MIPS KAOOS

Sequencing Projects SPP Consortium GSC MGI ESSA Genoscope

Arabidopsis thaliana genome view BLAST search Arabidopsis genome

Hits: Chr Match Map element Type Maps

V At5g44790 NM\_123847 GENE Gene

were provided to NCBI by TIGR on behalf of the collaborators.

**Map View: NM**

Arabidopsis thaliana Map View

Chromosome: I II III IV [ V ] Query: At5g44790 [clear]

Master Map: Gene

Total Genes On Chromosome: 6000 Region Displayed: 17,787K-17,796K bp Download/View Sequence/Evidence Genes Labeled: 1 Total Genes in Region

Region Shown: 17,787K Go

Map Viewer Help Arabidopsis Maps Help FTP

Data As Table View Maps&Options

ideogram master

Reverse Complement Strand View on plus strand Protein coding genes Hide Toolbar

Search for gene Find Refresh

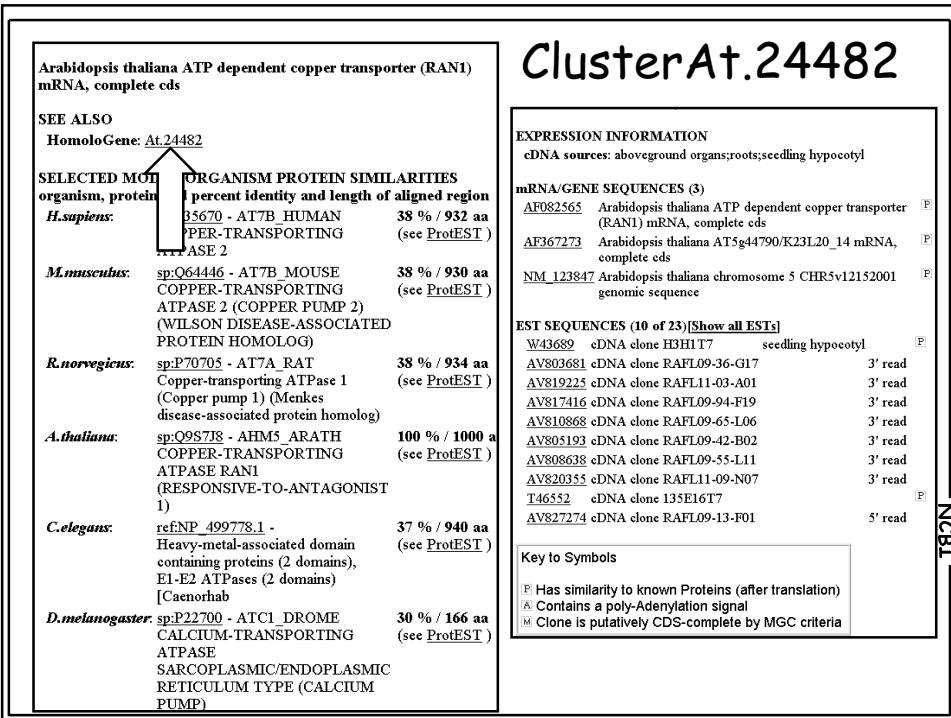
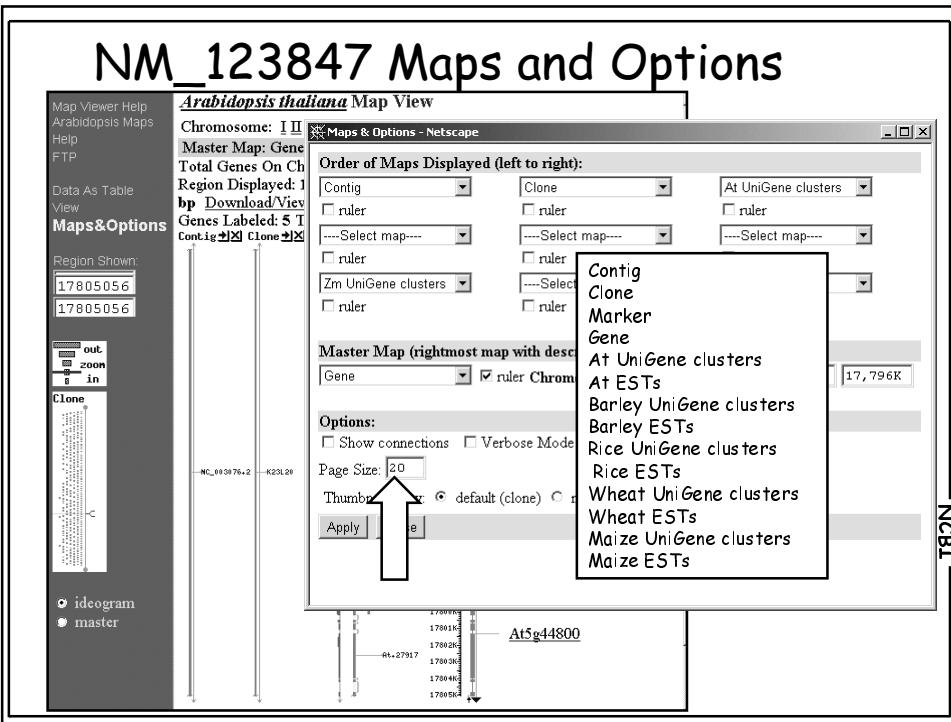
26689408 21M 19M 1 9M

Sequence:

```

1779910 TTGTTTTTTT TTTCCTGGAT TTTCCTTTCTT TCAAAATTTC CTCCTATGTT ATTACATTTA At5g44790 mRNA
1779950 CCAAAAGACA AAGCTCTTC AGTTTGTAA ATTCGGTTCTT TCAAGAGAC At5g44790 mRNA
1779980 CAAAGATA CTTCCCTCTT TATCTCTTCC GATTCTCAA CCTTCCCGTC TCTCTTCA At5g44790 mRNA
1779990 CGGCGAATTC CCCCGCTGTC ACCGGTTTA CCGCTTCTTC CGGTACCGG AGCTTCTTC At5g44790 mRNA
1779995 T C P S R R D L Q L T P V T G C S S
1779996 TCCGATTTA GTCATGTTA AGAGTTGGT CTTCCTGTT CTGATGTTA CGACGCCG At5g44790 mRNA
1779997 S O I S D H E E U G L L D S Y H N E R N
1779998 CGCGATGTTA TTTTGCTTA AATCCGAAGAA CGGAGGTTTC TTTCGGTTT AAGGAGTT At5g44790 mRNA
1779999 R D D I L T K I E E G R D V S G L R K I
17799999 CGGGTGGG TCGCCGGAT GAGTTGGT GCTTGTCTTA ATTCGGTTA AGCGCTTCA At5g44790 mRNA
177999999 Q V G C T G C H T T C A R C A S N S V E R A
1779999999 TGAAGGTTA ATGGCGCTT CGAACGCTT GTGCGTTGT TACGGATTC ACCGGATTC At5g44790 mRNA
17799999999 M H V N G U F K R S V A L L U N R A D V
177999999999 CTCCTGGCC CTTTTGTTA CGACGTTA TCTTGGTTA TTTTGGTT CTTATTTA At5g44790 mRNA
V F D P N L V K

```



**HomoloGene**

**Gene Info**

**CALCULATED ORTHOLOGS**

Listed below are UniGene clusters which are possible orthologs as determined by nucleotide sequence comparisons. The % ID below includes hyperlinks to the indicated alignments  
**MORE ▶**

Organism	Sequence	% ID	Sequence	Organism
Thale Cress	At.24482 AF082565	60.7	Rn.10554 U59245	Rat
Rat			↔ U59245	LocusLink
Zebrafish			24941 Atp7a	24941 Atp7a
Fly			Rn.10554 U59245	Human
Arabidopsis thaliana (thale cress)		89.5	Hs.606 Z94801	
Hordeum vulgare			↔ LocusLink	
			24941 Atp7a	538 ATP7A

A double headed arrow indicates that the pair ↔ represents a reciprocal best hit, the match is the best one for both organisms.

A single headed arrow indicates that the pair ← represents the best match for one of the organisms. The direction of the arrow indicates whether it is from or to the corresponding organism.

When present, red arrows point out a group of ▶ sequence matches which are part of a triplet, being consistent between more than two organisms.

**LocusLink**

**A single query interface to**

**UniGene**

**PubMed**

**HomoloGene**

**Links**

**Position**

**Xq13.2-q13.3 P O R G P H U V**

**Full report**

**538 Hs ATP7A ATPase, Cu++ transporting, alpha polypeptide (M**

**syndrome)**

**11977 Mm AtP7a ATPase, Cu++ transporting, alpha polypept Available for**

**24941 Rn AtP7a ATPase Hs human**

**transpor Mm mouse**

**polypep Rn rat**

**syndrom Dr zebrafish**

**11979 Mm AtP7b ATPase Dm fruit fly**

**transpor Dm HIV**

**polypep HIV**

**Map Viewer**

**OMIM**

**RefSeq**

**GenBank Accessions**

**R G P**

**dbSNP**

**0.0 cM P R G P H U**

# LocusLink ATP7A

Click to Display mRNA-Genomic Alignments (spanning 78777 bps)

PUB	OMIM	UNIGENE	MAP	VAR	HOMOL	GDR	HOMD	eL
UCSC	PROTEOME							

**Homo sapiens Official Gene Symbol and Name (HGNC)**  
ATP7A: ATPase, Cu++ transporting, alpha polypeptide (Menkes syndrome)  
LocusID: 538

**Overview**

**Protein Summary:** Copper transporting ATPase alpha polypeptide; important for copper efflux from cells; has six putative metal-binding motifs

**Locus Type:** gene with protein product, function known or inferred

**Product:** ATPase, Cu<sup>++</sup> transporting, alpha polypeptide

**Alternate Symbols:** MK, MNK, OHS

**Function:** Submit GeneRIF (All Pubs) ?

**Phenotype:**

- Cutis laxa, neonatal
- Menkes disease
- Occipital horn syndrome

**EC Number:** 3.6.1.36

**Gene Ontology™:**

Term	Evidence Source	Pub
Golgi apparatus	E	Proteome pm
copper ion transport	P	Proteome pm
copper-exporting ATPase	P	Proteome pm
integral plasma membrane protein	NR	Proteome

**Other Ontologies:**

Term	Evidence Source	Pub
Golgi	NR	Proteome pm
Integral membrane	NR	Proteome pm
Unspecified membrane	P	Proteome pm

**Relationships** ?

**Mouse Homology Maps:**

NCBI vs. MGD	X 44.00 cM	Atp7a	Hs	Mm
UCSC vs. MGD	X 44.00 cM	Atp7a	Hs	Mm

**Map Information**

Chromosome: X mv

Cytogenetic: Xq13.2-q13.3 HUGO

Markers:

Chr. X	CHLC.UTR_01549_L06133	mv
Chr. X	SGC31483	mv
Chr. -	AI023617	mv
Chr. -	GDB:435280	
Chr. -	GDB:435283	
Chr. -	GDB:435286	
Chr. -	GDB:435290	
Chr. -	GDB:435294	

NCBI Reference Sequences (RefSeq)

Category: PROVISIONAL

mRNA: NM\_000052

Protein: NP\_000043 ATPase, Cu<sup>++</sup> transporting, BL

alpha polypeptide

Domains: E1-E2 ATPase score: 532

Heavy-metal-associated domain score: 196

GenBank: L06133

Source: NCBI

Links:

**pm** PubMed

**mv** MapViewer

**sv** Sequence Viewer

**ev** Evidence Viewer

**BL** BLink

Category: NCBI Genome Annotation

Genomic: NT\_030881 sv mv ev

Config: Annotated transcripts/proteins for this locus:

Evidence: none

Model: XM\_013141

mRNA: XP\_013141

Protein: XP\_013141 BL

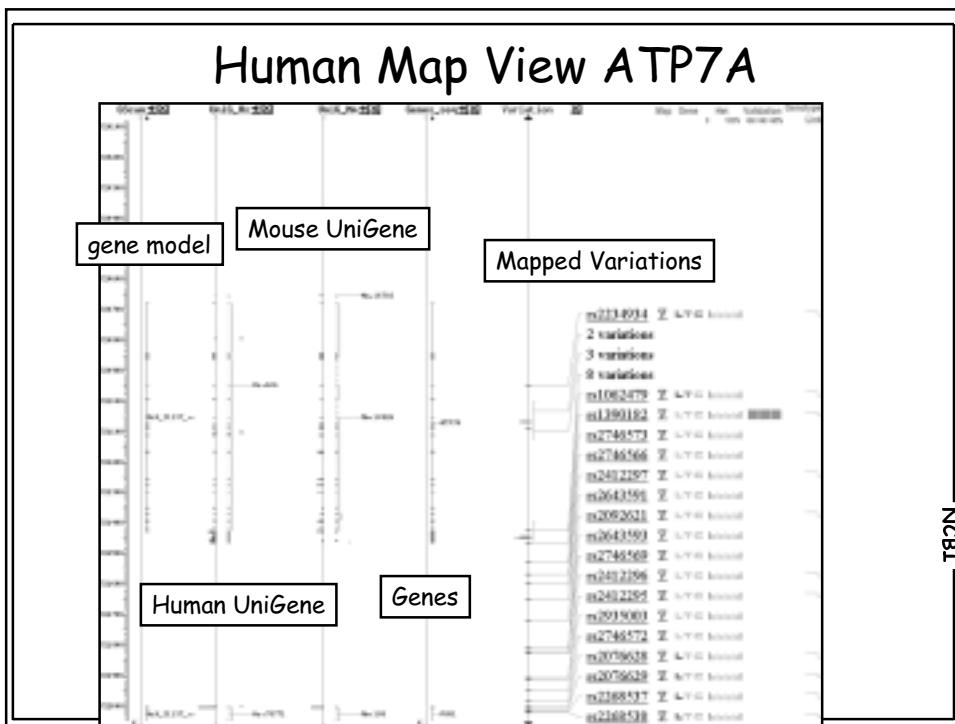
Domains: E1-E2 ATPase score: 589

Heavy-metal-associated domain score: 190

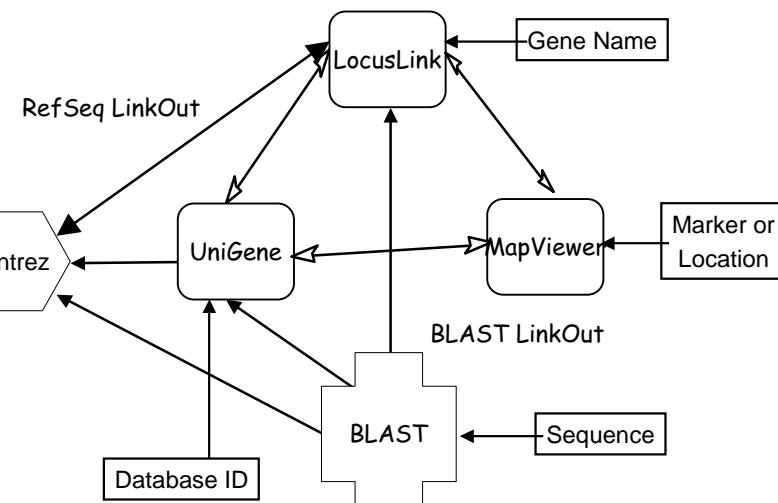
haloacid dehalogenase-like hydrolase score: 200

GenBank Sequences

Nucleotide Type Protein



## Genome Resources Integration



## Why do we need similarity searching?

- ◆ Identification and annotation
  - Incomplete or no annotations (GenBank)
  - Incorrectly annotated sequences
- ◆ Evolutionary relationships
  - homologous molecules **may** have similar functions

but it ain't necessarily so!

## Basic Local Alignment Search Tool

- Widely used similarity search tool
- Heuristic approach based on Smith Waterman algorithm
- Finds best local alignments
- Provides statistical significance
- All combinations (DNA/Protein) query and database.
  - DNA vs DNA
  - DNA translation vs Protein
  - Protein vs Protein
  - Protein vs DNA translation
  - DNA translation vs DNA translation
- [www](#), email server, standalone, and network clients

NCBI

## How BLAST Works

- Make lookup table (hash table) for query
- Scan database for hits
- Ungapped extensions of hits
- Gapped extensions (no traceback)
- Gapped extensions (traceback)

NCBI

## Look Up Table (Hash Table)

Query: **GTQITVEDLFYNIATRRKALKN**

**GTQ**

Word Size = 3 Adjustable

**TQI**

2 or 3 for protein ( 3 default)

**QIT**

> 7 for blastn searches ( 11 default)

**Neighborhood Words**

**ITV** -> **LTV, MTV, ISV, LSV, MSV**

**TVE** **IAV, LAV, MAV, ITL, etc.**

**VED**

**EDL**

**DLF**

**LFY**

**FYN**

Make table  
for both query  
and database

NCBI

## Messy Details

ATCGCCATGCTTAATTGGGCTT

**CATGCTTAATT** exact word match

one hit

- Nucleotide BLAST looks for exact matches
- Protein BLAST requires two hits

**GTQITVEDLFYNI**

**SEI**

**YYN**

neighborhood words

**{ }**      **{ }**

two hits

NCBI

## More Details (BLAST options)

```
-W Word size
-f Threshold for extending hits
-X X dropoff value for gapped alignment (in bits)
-y Dropoff (X) for blast extensions in bits
-Z X dropoff value for final gapped alignment (in bits)
-A Multiple Hits window size (zero for single hit algorithm)
-e Expectation value (E) default = 10.0
-q Penalty for a nucleotide mismatch (blastn only) default = -3
-r Reward for a nucleotide match (blastn only) default = 1
-v Number of database one-line descriptions
-b Number of database alignments
```

NCBI

## An alignment that BLAST can't find

```
1 GAATATATGAAGACCAAGATTGCAGTCCTGCTGGCCTGAACCACGCTATTCTTGCTGTTG
   ||| ||| |||| | ||| ||| ||| | ||| | |||| |||| | ||| | ||| |
1 GAGTGTACGTGAGCCCGAGTGTAGCAGTGAAGATCTGGACCACGGTGTACTCGTTGTCG

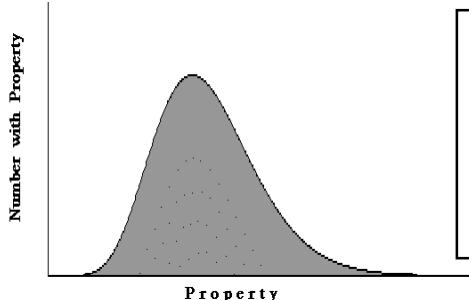
61 GTTACGGAACCGAGAATGGTAAAGACTACTGGATCATTAAGAACCTCCTGGGAGCCAGTT
   ||| ||| ||| | ||| ||| | ||| ||| | ||| ||| | ||| ||| | | |
61 GCTATGGTGTAAAGGTGGAAAGAAGTACTGGCTCGTCAAGAACAGCTGGCTGAATCCT

121 GGGGTGAACAAGGTTATTCAGGCTTGCTCGTGGTAAAAAC
   |||| | ||| | | | | ||| | ||| | ||| |
121 GGGGAGACCAAGGCTACATCCTATGTCCCGTACAACAAAC
```

NCBI

## Local Alignment Statistics

High scores of local alignments between two random sequences follow Extreme Value Distribution



For ungapped alignments:

Expected number with score S or greater

$$E = Kmne^{-\lambda S}$$

or

$$E = mn2^{-S'}$$

K = scale for search space

$\lambda$  = scale for scoring system

$$S' = \text{bitscore} = (\lambda S - \ln K) / \ln 2$$

<http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html>

NCBI

## Scoring Systems

• Nucleic acids identity matrix

• Proteins

• Position Independent Matrices

• PAM Matrices (Percent Accepted Mutation)

• Implicit model of evolution

• Higher PAM number all calculated from PAM1

• PAM250 widely used

• BLOSUM Matrices (BLOck SUbstitution Matrices)

• Empirically determined from alignment of conserved blocks

• Each includes information up to a certain level of identity

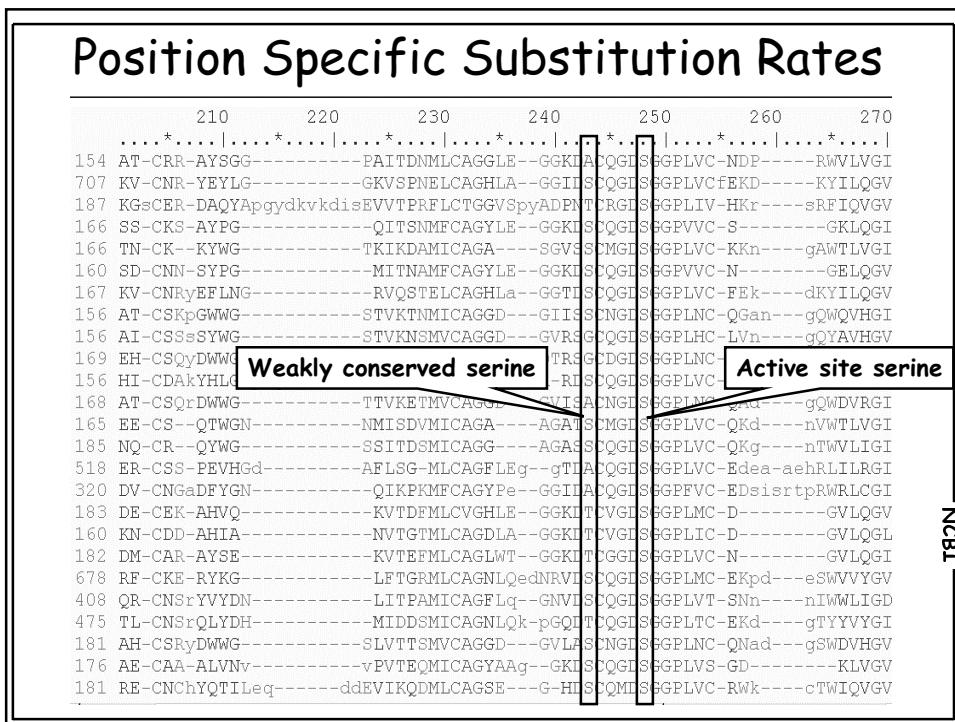
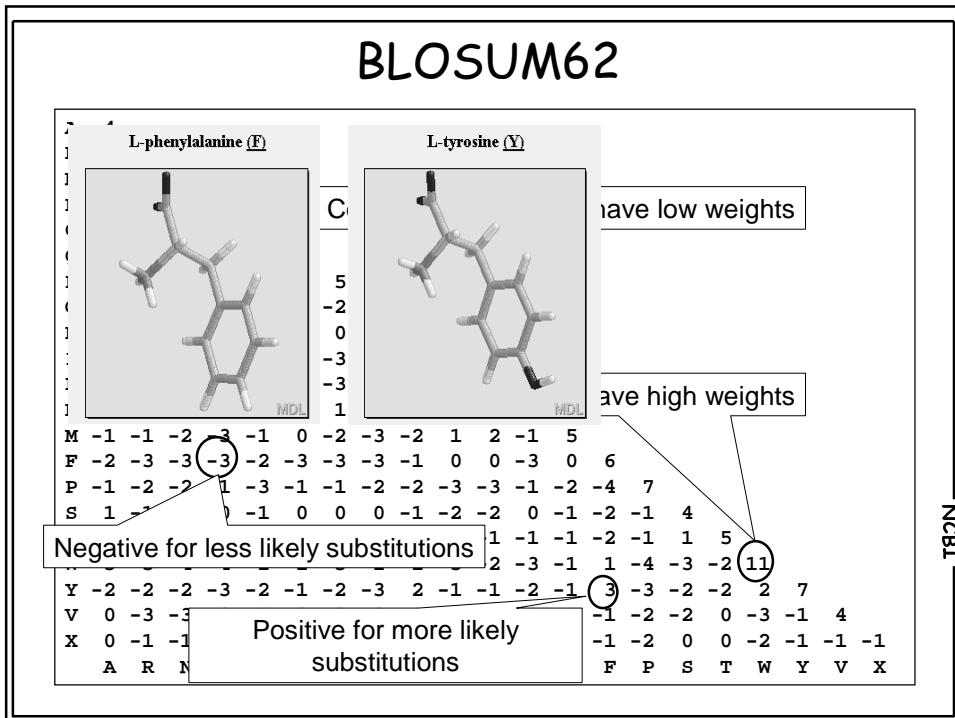
• BLOSUM62 widely used

• Position Specific Score Matrices (PSSMs)

• PSI and RPS BLAST

	A	G	C	T
A	+1	-3	-3	-3
G	-3	+1	-3	-3
C	-3	-3	+1	-3
T	-3	-3	-3	+1

NCBI



## Position Specific Score Matrix (PSSM)

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V		
206	D	0	-2	0	2	-4	2	4	-4	-3	-5	-4	0	-2	-6	1	0	-1	-6	-4	-1	
207	G	-2	-1	0	-2	-4	-3	-3	6	-4	-5	-5	0	-2	-3	-2	-2	-1	0	-6	-5	
208	V	-1	1	-3	-3	-5	-1	-2	6	-1	-4	-5	1	-5	-6	-4	0	-2	-6	-4	-2	
209	I	-3	3	-3	-4	-6	0	-1	-4	-1	2	-4	6	-2	-5	-5	-3	0	-1	-4	0	
210	S	-2	-5	0	8	-5	-3	-2	-1	-4	-7	-6	-4	-6	-7	-5	1	-3	-7	-5	-6	
211	S	4	-4	-4	-4	-4	-1	-4	-2	-3	-3	-5	-4	-4	-5	-1	4	3	-6	-5	-3	
212	C	-4	-7	-6	-7	12	Serine scored differently in these two positions										-7	-4	-4	-5	0	-4
213	N	-2	0	2	-1	-6											-5	-1	-3	-3	-4	-3
214	G	-2	-3	-3	-4	-4											-6	-3	-5	-6	-6	-6
215	D	-5	-5	-2	9	-7	-4	-1	-5	-5	-7	-7	-4	-7	-7	-5	-4	-4	-8	-7	-7	
216	S	-2	-4	-2	-4	-4	-3	-3	-3	-4	-6	-6	-3	-5	-6	-4	7	-2	-6	-5	-5	
217	G	4	-6	-4	-5	-6	-5	-6	8	-6	-8	-7	-5	-6	-7	-6	-4	-5	-6	-7	-7	
218	G	-2	-2	-2	-2	-2	-2	-2	6	-7	-7	-5	-6	-7	-6	-2	-4	-6	-7	-7	-7	
219	P	-2	-2	-2	-2	-2	-2	-2	6	-6	-7	-4	-6	-7	9	-4	-4	-7	-7	-6	-6	
220	L	-4	-6	-7	-7	-5	-5	-6	-7	0	-1	6	-6	1	0	-6	-6	-5	-5	-4	0	
221	N	-1	-6	0	-6	-4	-4	-6	-6	-1	3	0	-5	4	-3	-6	-2	-1	-6	-1	6	
222	C	0	-4	-5	-5	10	-2	-5	-5	1	-1	-1	-5	0	-1	-4	-1	0	-5	0	0	
223	Q	0	1	4	2	-5	2	0	0	0	-4	-2	1	0	0	0	-1	-1	-3	-3	-4	
224	A	-1	-1	1	3	-4	-1	1	4	-3	-4	-3	-1	-2	-2	-3	0	-2	-2	-2	-3	

NCBI

## Gapped Alignments

- Gapping provides more biologically realistic alignments
- Statistical behavior not completely understood for gapped alignments
- Gapped BLAST parameters must be found by simulations for each matrix
- Affine gap costs =  $-(a+bk)$   
 $a$  = gap open penalty     $b$  = gap extend penalty  
A gap of length 1 receives the score  $-(a+b)$

NCBI

## Scores

	V	D	S	-	C	Y	
	V	E	T	L	C	F	
BLOSUM62	+4	+2	+1	-12	+9	+3	<u>7</u>
PAM30	+7	+2	0	-10	+10	+2	<u>11</u>

NCBI

WWW BLAST

NCBI

**NCBI** PubMed Entrez BLAST OMIM Taxonomy Structure

NCBI SITE MAP BLAST info BLAST overview Frequently Asked Questions New/Noteworthy Receive e-mail with BLAST announcements BLAST course BLAST tutorial BLAST references URL API documentation HTML format PDF format PostScript format FTP BLAST FTP site Credits BLAST Credits

## BLAST

**What's NEW in BLAST®**

September 26, 2001: New pages are now available for Mouse, Rat and Fugu genomes in the "Genomic BLAST pages" section.

**Results from the Protein Sequence Information Survey**

NCBI's Protein Sequence Information Survey Results are here. Thank you for participating.

**Nucleotide BLAST**

- Standard nucleotide-nucleotide BLAST [blastn]
- MEGABLAST
- Search for short nearly exact matches

**Protein BLAST**

- Standard protein-protein BLAST [blastp]
- PSI- and PHI-BLAST
- Search for short nearly exact matches

**Translated BLAST Searches**

- Nucleotide query - Protein db [blastx]
- Protein query - Translated db [tblastn]
- Nucleotide query - Translated db [tblastx]

**Search for conserved domains**

- Search the Conserved Domain Database using RPS-BLAST
- Search by domain architecture [DART]

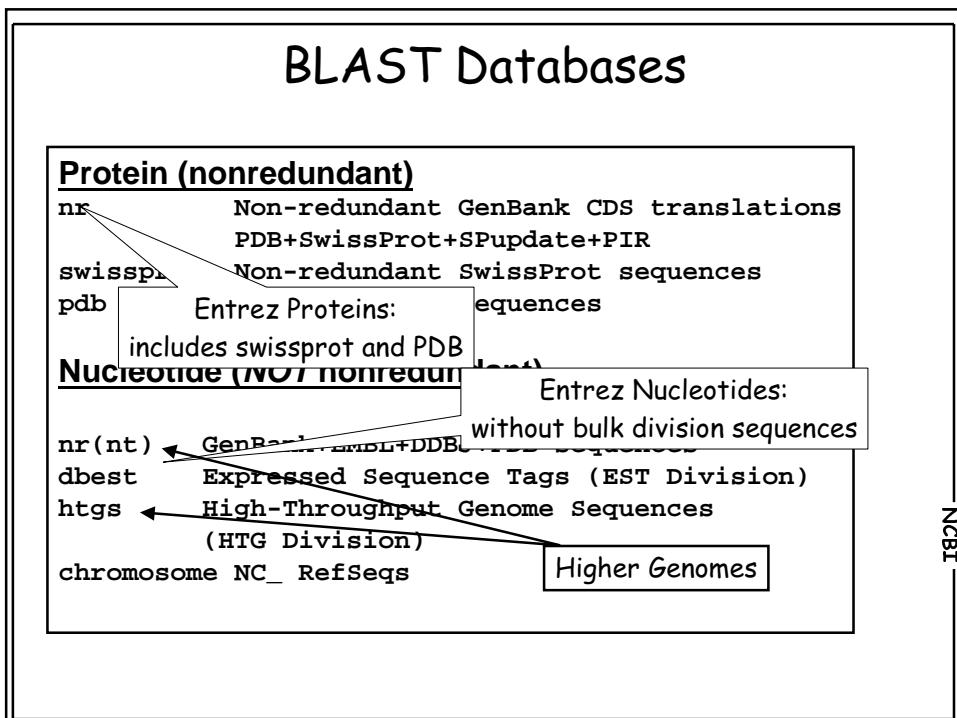
**Pairwise BLAST**

- BLAST 2 Sequences

**Genomic BLAST pages**

- Human Genome
- Microbial Genomes
- Arabidopsis thaliana
- Other eukaryotes
- Mouse Genome
- Rat Genome
- Fugu rubripes

**Web BLAST**



## Protein BLAST Page

NCBI

**protein-protein BLAST**

Nucleotide    Protein    Translations    Retrieve results for an RID

Search >Mutated in Colon Cancer  
 IETVYAAALPKNTHPFLYLSLEISPQNVNDVNVPDKHEVHFLHEESILER  
 VQQHIESKLLGSNSRMYFTQTLLPGLAGPSGEMVKSTTSLTSSSTSGSS  
 DVKVAHQMVRTDSREQKLDAFLQPLSKPLSS

Set subsequence From: [ ] To: [ ]

Choose database swissprot    **Protein database**

Do CD-Search

Now: **BLAST!** or **Reset query** **Reset all**

## BLAST Formatting Page

NCBI

**formatting BLAST**

Nucleotide    Protein    Translations    Retrieve results for an RID

Your request has been successfully submitted and put into the Blast Queue.

Query = Mutated in Colon Cancer (131 letters)

---

Putative conserved domains have been detected

Click on the image below for detailed CD-Search results

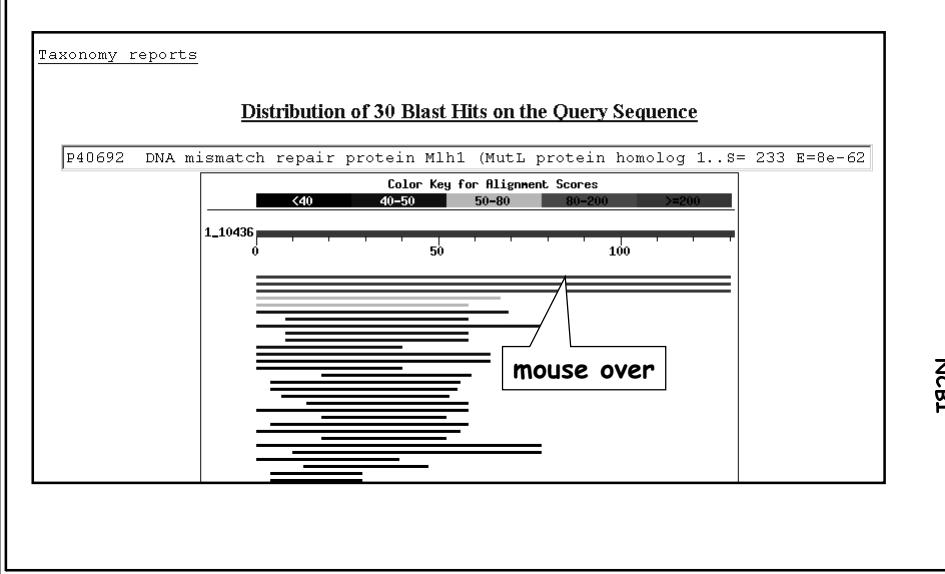


The request ID is 989277607-10255-30769

**Format!** or **Reset all**

The results are estimated to be ready in 4 seconds but may be done sooner.

## BLAST Output: Graphic



BLAST Output: Descriptions

sorted by e values

Sequences producing significant alignments

	Score	E	(bits)	Value
gi 730028 sp P40692 MLH1 HUMAN DNA mismatch repair protein ...	233	8e-62		
gi 13878583 sp Q9JK91 MLH1 MOUSE DNA m...	214	4e-56		
gi 13878571 sp P97679 MLH1 RAT DNA mis...	212	1e-55		
gi 17090561 sp P38920 MLH1 YEAST MUTL protein homolog 1 (DNA...	72	2e-13		
gi 11710801 sp P44494 MUTL HAEIN DNA mismatch repair protein...	54	7e-08		
gi 13431695 sp P57886 MUTL PASMU DNA mismatch repair protei...	42	1e-06		
gi 117090561 sp P44494 MUTL THEMA DNA mismatch repair protein...	48	4e-06		
gi 117090561 sp P44494 MUTL BACHD DNA mismatch repair protein...	46	1e-05		
gi 117090561 sp P44494 MUTL ECOLI DNA mismatch repair protein...	44	5e-05		
gi 127553 sp P14161 MUTL SALTY DNA mismatch r...	44	7e-05		
gi 6225738 sp Q9ZC88 MUTL RICPR DNA mismatch r...	40	7e-04		
gi 14194944 sp Q9PJG5 MUTL CHLMU DNA mismatch r...	40	0.001		
gi 18928218 sp O84579 MUTL CHLTR DNA mismatch repair protein...	39	0.001		
gi 20043258 sp Q9KV13 MUTL VIBCH DNA mismatch repair protein...	39	0.002		
gi 13631230 sp Q9RP66 MUTL CAUCR DNA mismatch repair protein...	39	0.002		
gi 18928214 sp O51229 MUTL BORBU DNA mismatch repair protein...	39	0.002		
gi 1709188 sp P49850 MUTL BACSU DNA mismatch repair protein...	38	0.005		
gi 8039787 sp O83325 MUTL TREPA DNA mismatch repair protein...	36	0.013		
gi 19856116 sp P14160 HEX DNA mismatch repair protein...	36	0.020		
gi 3914082 sp P70754 MUTL DNA mismatch repair protein...	35	0.020		
gi 11386926 sp P57633 MUTL DNA mismatch repair protein...	35	0.026		
gi 18928240 sp Q9Z794 MUTL CHLPN DNA mismatch repair protein...	34	0.035		
gi 1709684 sp P54280 PMS1 SCHPO DNA mismatch repair protein...	33	0.16		
gi 13914001 sp Q67510 MUTL AQUAE DNA mismatch repair protein...	33	0.16		
gi 1709685 sp P54278 PMS2 HUMAN PMS1 protein homolog 2 (DNA...	33	0.16		
gi 1709686 sp P54279 PMS2 MOUSE PMS1 PROTEIN HOMOLOG 2 (DNA...	32	0.24		
gi 18928222 sp P73349 MUTL SYNY3 DNA mismatch repair protein...	31	0.60		
gi 1709683 sp P54277 PMS1 HUMAN PMS1 protein homolog 1 (DNA...	30	0.85		
gi 126232 sp P02239 LGB1 LUPLU Leghemoglobin I	30	1.2		
gi 126238 sp P02240 LGB2 LUPLU Leghemoglobin II	28	4.1		

Default e value cutoff 10

Bacterial mismatch repair proteins

## TaxBLAST: Taxonomy Reports

<u><a href="#">Haemophilus influenzae</a></u> [g-proteobacteria] taxid 927 gi 1171010 sp P46494 MUTL_HASIN DNA mismatch repair protei...	<u>52</u>	<u>7e-08</u>
<u><a href="#">Pastorevella maltocida</a></u> [g-proteobacteria] taxid 747 gi 113431675 sp P57864 MUTL_DASHU DNA mismatch repair protei...	<u>49</u>	<u>1e-06</u>
<u><a href="#">Thermotoga maritima</a></u> [thermotogaes] taxid 2336 gi 10920224 sp P74825 MUTL_THEME DNA mismatch repair protei...	<u>47</u>	<u>4e-06</u>
<u><a href="#">Bacillus halodurans</a></u> [low GC Gram+] taxid 86665 gi 120139217 sp Q9AC11 MUTL_BACHD DNA mismatch repair protei...	<u>46</u>	<u>1e-05</u>
<u><a href="#">Escherichia coli</a></u> [enterobacteria] taxid 562 gi 127552 sp P23367 MUTL_ECOLI DNA mismatch repair protein...	<u>44</u>	<u>5e-05</u>
<u><a href="#">Salmonella typhimurium</a></u> [enterobacteria] taxid 602 gi 127553 sp P14161 MUTL_SALTY DNA mismatch repair protein...	<u>43</u>	<u>7e-05</u>
<u><a href="#">Rickettsia prowazekii</a></u> [a-proteobacteria] taxid 782 gi d225738 sp Q92C88 MUTL_RICPR DNA mismatch repair protei...	<u>40</u>	<u>7e-04</u>
<u><a href="#">Chlamydia muridarum</a></u> (agent of mouse pneumonitis) [chlamydias] taxid 83563 gi 14194944 sp Q9RJG5 MUTL_CMLMU DNA mismatch repair protei...	<u>40</u>	<u>0.001</u>
<u><a href="#">Chlamydia trachomatis</a></u> [chlamydias] taxid 813 gi 8926218 sp Q84579 MUTL_CHLTR DNA mismatch repair protei...	<u>39</u>	<u>0.001</u>
<u><a href="#">Vibrio cholerae</a></u> [g-proteobacteria] taxid 666 gi 20043250 sp Q9EV13 MUTL_VIBCH DNA mismatch repair protei...	<u>38</u>	<u>0.002</u>
<u><a href="#">Caulobacter vibrioides</a></u> [a-proteobacteria] taxid 155892 gi 134331230 sp Q9F766 MUTL_CAOCH DNA mismatch repair protei...	<u>38</u>	<u>0.002</u>
<u><a href="#">Borrelia burgdorferi</a></u> (lyme disease spirochete) [spirochetes] taxid 139 gi 8926214 sp Q51329 MUTL_BORSP DNA mismatch repair protei...	<u>38</u>	<u>0.002</u>

NCBI

## BLAST Output: Alignments

```
>gi|127552|sp|P23367|MUTL_ECOLI  DNA mismatch repair protein mutL
Length = 615

Score = 44.3 bits (103), Expect = 5e-05
Identities = 25/59 (42%), Positives = 33/59 (55%), Gaps = 8/59 (13%)

Query: 9   LPKNTHPFLYLSLEISPQNVDVNHPPTKHEVHF-----LHE---ESILERVQQHIESKL 59
          L + P   L LEI P   VDVNVHP KHEV F      +H+    + +L   +QQ +E+ L
Sbjct: 280 LGADQQPAFVLYLEIDPHQVDVNHPAKHEVRFHQSRLVHDFIYQGVLSVLQQQLETPL 338
```

NCBI

## BLAST Output: Alignments

```

>gi|730028|sp|P40692|MLH1_HUMAN DNA mismatch repair protein Mlh1 1)
Length = 756

Score = 233 bits (593), Expect = 8e-62
Identities = 117/131 (89%), Positives = 117/131 (89%)

Query: 1 IETVYAAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL 60
        IETVYAAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL
Sbjct: 276 IETVYAAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL 335

Query: 61 GSNSSRMYFTQTLLPGLAGPSGEMVKXXXXXXXXXXXXXXDKVYAHQMVRTDSREQKLDA 120
        GSNSSRMYFTQTLLPGLAGPSGEMVK                               DKVYAHQMVRTDSREQKLDA
Sbjct: 336 GSNSSRMYFTQTLLPGLAGPSGEMVKSTTSLTSSSTSGSSDKVYAHQMVRTDSREQKLDA 395

Query: 121 FLQPLSKPLSS 131 low complexity sequence filtered
        FLQPLSKPLSS
Sbjct: 396 FLQPLSKPLSS 406

```

NCBI

## Results from nr

Sequences producing significant alignments:	(bits)	Value
gi 604369 gb AAA85687.1  (U17857) hMLH1 gene product [Homo ...	233	3e-61
>gi 4557757 ref NP_000240.1  (NM_000249) mutL homolog 1; mutL (E. coli) homolog 1; colli) homolog 1 (colon cancer, nonpolyposis type 2) [Homo sapiens]		
gi 730028 sp P40692 MLH1_HUMAN DNA mismatch repair protein Mlh1 (MutL protein homolog 1)		
gi 631299 pir  S43085 DNA mismatch repair protein MLH1 - human		
gi 463989 gb AAC50285.1 (U07343) hMLH1 [Homo sapiens]		
gi 1079787 gb AAA82079.1 (U40978) DNA mismatch repair protein homolog [Homo sapiens]		
gi 13905126 gb AAH06850.1 AAH06850 (BC006850) mutL (E. coli) homolog 1 type 2 [Homo sapiens]		
gi 741682 prf  2007430A DNA mismatch repair protein [Homo sapiens]		
Length = 756		
Score = 233 bits (593), Expect = 4e-61		
Identities = 117/131 (89%), Positives = 117/131 (89%)		
Query: 1 IETVYAAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL 60		
IETVYAAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL		
Sbjct: 276 IETVYAAAYLPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL 335		
gi 10272051 ref NP_438240.1  (NC_000907) DNA mismatch repair... 54 3e-07		
gi 19173567 ref NP_597370.1  (NC_003232) DNA MISMATCH REPAIR... 52 9e-07		
gi 13543339 gb AAH05833.1 AAH05833 (BC005833) Similar to mu... 50 5e-06		
gi 15602769 ref NP_245841.1  (NC_002663) MutL [Pasteurella ... 50 6e-06		
gi 15642797 ref NP_227838.1  (NC_000853) DNA mismatch repai... 48 2e-05		

NCBI

## tblastn Results Against ESTs

```

>gi|12794555|emb|AL531062.1|AL531062 AL531062 LTI_NFL001_NBC4 Homo sapiens
  cDNA clone CS0DM005YM23 5
    prime.
    Length = 878

  Score = 167 bits (422), Expect(3) = 1e-42
  Identities = 81/82 (98%), Positives = 81/82 (98%)
  Frame = +2

  Query: 1 IETVYAAYLKPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL 60
          IETVYAAYLKPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL
  Sbjct: 512 IETVYAAYLKPKNTHPFLYLSLEISPQNVDVNVHPTKHEVHFLHEESILERVQQHIESKLL 691

  Query: 61 GSNSSRMYFTQTLLPGLAGPSG 82
          GSNSSRMYFTQTLLPGLAGP G
  Sbjct: 692 GSNSSRMYFTQTLLPGLAGPLG 757

  Score = 24.3 bits (51), Expect(3) = 1e-42
  Identities = 11/26 (42%), Positives = 11/26 (42%)
  Frame = +1

  Query: 80 PSGEMVKXXXXXXXXXXXXXDKVVA 105
          PSG MVK DVKVVA
  Sbjct: 748 PSG*MVKSTTSLTSSSTSGSSDKVVA 825

```

combined expect for  
hits to multiple frames

NCBI

## Results against PDB - Finding a model template

MMDB Id: 9958 PDB Id: Cn3D 3.0

Protein Chains:	B, A
MEDLINE:	<a href="#">PubMed</a>
Taxonomy:	B, A <a href="#">Escherichia coli</a>
PDB Authors:	Y. Wei
PDB Deposition:	11-Jan-99
PDB Class:	Dna Mismatch
PDB Title:	MutL Complex
Sequence Neighbors:	<a href="#">B, A</a>
Structure Neighbors:	<a href="#">B, A, A.1, A</a>

View / Save Structure      NEW

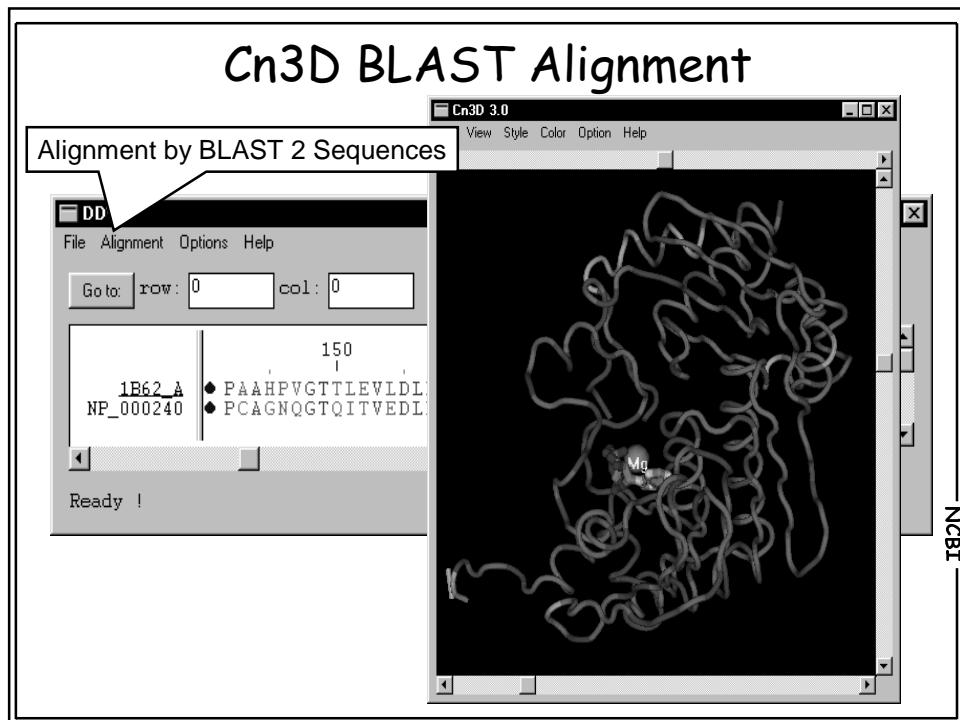
Options:      Viewer:

- Launch View
- See File
- Save File
- Cn3D (asn.1)
- Cn3D v1.0 (a)
- Mage
- RasMol (PDB)

Value

1e-05  
4e-05  
2.0

NCBI



# PSI-BLAST

Confirming relationships of purine nucleotide metabolism proteins

NCBI

**PSI BLAST**

**Search**

```
>gi|113340|sp|P03958|ADA MOUSE ADENOSINE DEAMINASE (ADENOSINE AMI...  
MAQTPAEI  
VIAGCRE  
EQAFGKIE  
RTVHAGE  
VRFKNDR
```

Set subsequence From:

Choose database **swissprot**

Do CD-Search

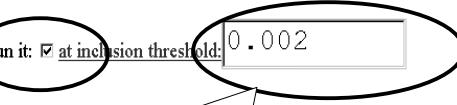
Now: **BLAST!** 

**Format**

Show  Graphical Overview  NCBI-gi Alignment  in **HTML**  format

Number of: Descriptions **500** Alignments **250**

Alignment view **Pairwise**

PSI-BLAST Run it:  at inclusion threshold: **0.002** 

Autoformat  **e value cutoff for PSSM**



**PSI RESULTS: Initial BLAST Run**

Sequences with E-value BETTER than threshold

	Score	E
	(bits)	Value
Sequences producing significant alignments:		
NEW <input checked="" type="checkbox"/> <a href="#">gi 113340 sp P03958 ADA MOUSE</a> ADENOSINE DEAMINASE (ADENOSINE AMI...	679	0.0
NEW <input checked="" type="checkbox"/> <a href="#">gi 5902736 sp P56658 ADA BOVIN</a> ADENOSINE DEAMINASE (ADENOSINE AMI...	608	e-174
NEW <input checked="" type="checkbox"/> <a href="#">gi 113339 sp P00813 ADA HUMAN</a> ADENOSINE DEAMINASE (ADENOSINE AMI...	600	e-171
NEW <input checked="" type="checkbox"/> <a href="#">gi 2506342 sp P22333 ADD ECOLI</a> ADENOSINE DEAMINASE (ADENOSINE AMI...	151	2e-36
NEW <input checked="" type="checkbox"/> <a href="#">gi 1703166 sp P53909 ADA YEAST</a> ADENOSINE DEAMINASE (ADENOSINE AMI...	103	7e-22
NEW <input checked="" type="checkbox"/> <a href="#">gi 1703170 sp P53984 ADD STRVG</a> ADENOSINE DEAMINASE (ADENOSINE AMI...	101	2e-21

Run PSI-Blast iteration 2 

Sequences with E-value WORSE than threshold

<input type="checkbox"/> <a href="#">gi 250619 sp P09030 EX3 ECOLI</a> EXODEOXYRIBONUCLEASE III (EXONUCL...	30	9.3
---	----	-----

Run PSI-Blast iteration 2



**First PSSM Search**

Sequences with E-value BETTER than threshold			
Sequences producing significant alignments:		Score	E (bits) Value
<input checked="" type="checkbox"/>	gi 5902736 sp P56658 ADA_BOVIN ADENOSINE DEAMINASE (ADENOSINE AM... 618 e-177		
<input checked="" type="checkbox"/>	gi 113339 sp P00813 ADA_HUMAN ADENOSINE DEAMINASE (ADENOSINE AMI... 615 e-176		
<b>Other purine nucleotide metabolizing enzymes not found by ordinary BLAST</b>			
<input checked="" type="checkbox"/>	gi 2506342 sp P22333 ADD_ECOLI ADENOSINE DEAMINASE 3 (AMP DEAMINASE IS... 452 e-127		
<input checked="" type="checkbox"/>	gi 1703166 sp P53909 ADA_YEAST ADENOSINE DEAMINASE 3 (AMP DEAMINASE... 438 e-123		
<input checked="" type="checkbox"/>	gi 1703170 sp P53984 ADD_STRVG ADENOSINE DEAMINASE 3 (AMP DEAMINASE AM... 419 e-117		
NEW	gi 399033 sp Q01432 AMD3_HUMAN AMP DEAMINASE 3 (AMP DEAMINASE IS... 53 1e-06		
NEW	gi 1351916 sp P15274 AMDm_YEAST AMP DEAMINASE (MYOADENYLATE DEAM... 51 4e-06		
NEW	gi 2494044 sp O09178 AMD3_RAT AMP DEAMINASE 3 (AMP DEAMINASE ISO... 49 1e-05		
NEW	gi 2494043 sp O08739 AMD3_MOUSE AMP DEAMINASE 3 (AMP DEAMINASE I... 49 2e-05		
NEW	gi 113698 sp P10759 AMD1_RAT AMP DEAMINASE 1 (MYOADENYLATE DEAMI... 46 8e-05		
NEW	gi 1703262 sp P50998 AMDm_SCHPO AMP DEAMINASE (MYOADENYLATE DEAM... 46 9e-05		
NEW	gi 399031 sp Q01433 AMD2_HUMAN AMP DEAMINASE 2 (AMP DEAMINASE IS... 46 1e-04		
NEW	gi 113697 sp P23109 AMD1_HUMAN AMP DEAMINASE 1 (MYOADENYLATE DEA... 44 4e-04		
Run PSI-Blast iteration 3			

NCBI

**Third PSSM Search: Convergence**

Sequences with E-value BETTER than threshold			
Sequences producing significant alignments:		Score	E (bits) Value
<input checked="" type="checkbox"/>	gi 5902736 sp P56658 ADA_BOVIN ADENOSINE DEAMINASE (ADENOSINE AM... 581 e-166		
<input checked="" type="checkbox"/>	gi 113339 sp P00813 ADA_HUMAN ADENOSINE DEAMINASE (ADENOSINE AMI... 578 e-165		
<input checked="" type="checkbox"/>	gi 113340 sp P03958 ADA_MOUSE ADENOSINE DEAMINASE (ADENOSINE AMI... 569 e-162		
<input checked="" type="checkbox"/>	gi 2506342 sp P22333 ADD_ECOLI ADENOSINE DEAMINASE 3 (AMP DEAMINASE IS... 569 e-162		
<input checked="" type="checkbox"/>	gi 1703166 sp P53909 ADA_YEAST ADENOSINE DEAMINASE 3 (AMP DEAMINASE... 569 e-162		
<input checked="" type="checkbox"/>	gi 1703170 sp P53984 ADD_STRVG ADENOSINE DEAMINASE 3 (AMP DEAMINASE... 569 e-162		
<input checked="" type="checkbox"/>	gi 399031 sp Q01433 AMD2_HUMAN AMP DEAMINASE 2 (AMP DEAMINASE IS... 569 e-162		
<input checked="" type="checkbox"/>	gi 1703262 sp P50998 AMDm_SCHPO AMP DEAMINASE (MYOADENYLATE DEAM... 569 e-162		
<input checked="" type="checkbox"/>	gi 2494044 sp O09178 AMD3_RAT AMP DEAMINASE 3 (AMP DEAMINASE ISO... 569 e-162		
<input checked="" type="checkbox"/>	gi 1399033 sp Q01432 AMD3_HUMAN AMP DEAMINASE 3 (AMP DEAMINASE IS... 568 e-162		
<input checked="" type="checkbox"/>	gi 2494043 sp O08739 AMD3_MOUSE AMP DEAMINASE 3 (AMP DEAMINASE I... 568 e-162		
<input checked="" type="checkbox"/>	gi 113698 sp P10759 AMD1_RAT AMP DEAMINASE 1 (MYOADENYLATE DEAMI... 568 e-162		
<b>Just below threshold, another nucleotide metabolism enzyme</b>			
<input type="checkbox"/>	sp P25524 CODA_ECOLI CYTOSINE DEAMINASE (CYTOSINE AMINOHYDROLASE) 38 0.041		
<input checked="" type="checkbox"/>	gi 113697 sp P23109 AMD1_HUMAN AMP DEAMINASE 1 (MYOADENYLATE DEAM... 208 2e-33		
<input checked="" type="checkbox"/>	gi 399032 sp Q02356 AMD2_RAT AMP DEAMINASE 2 (AMP DEAMINASE ISO... 134 2e-31		
<input checked="" type="checkbox"/>	gi 586396 sp P38150 YB92_YEAST HYPOTHETICAL 92.9 KD PROTEIN IN S... 131 1e-30		
<input checked="" type="checkbox"/>	gi 731934 sp P40361 YJHO_YEAST HYPOTHETICAL 104.3 KD PROTEIN IN ... 129 6e-30		

NCBI

**PHI BLAST**

Search

```
>gi|231729|sp|P30429|CED4_CAEEL CELL DEATH PROTEIN 4
MLCEIECRALSTAHTRLIHDFEPRALTYLEGKNIFTEDHSELISKMSTRLERIANFLRIYRRQASE
LIDFFNYNNQSHLADFLLEDYIDFAINEPDLLRPVVIAPQFSRQMLDRKLLLGNVPKQMTCYIREYHV
IKKLDEMCDDLSFLLFHGRAGSGKSVIASQALSKSDQLIGINYDSIWWLKDSGTAPKSTFDLFTDI
LKSEDLLNFPSPVEHTSVVLKRMICNALIDRPNTLFVFDDVVQETIRWAQELRLRCLVTTRDVEI
ASQTCFIEVTSLEIDECYDFLEAYGMPPVGKEKEEDVLNKTELSSGNPATLMMFFKSCEPKTFEK
```

Set subsequence From:  To:

Phi pattern **[GA]xxxxGK[ST]**

Do CD-Search

Now: **BLAST!** or

NCBI

**Conserved Domain Search**

**CDD**

KMICKHKNIINLLGACTOD...VIVIV...SCHCNLREVIQARRPPGLKSYMMWSHNDP...137  
TQL-EHSNLVOLLGVIV...SLW...100-GR...150...3  
TQL-EHSNLVOLLGVIV...GIV...100-GR...150...3  
KGFA...V...WSP...150-GR...200...52  
QEVSHPNVVIKLLGACTS...EXPLILITARY-SLRV...KIL...240...52

**NCBI**      **CD-Search**      **Entrez** ?

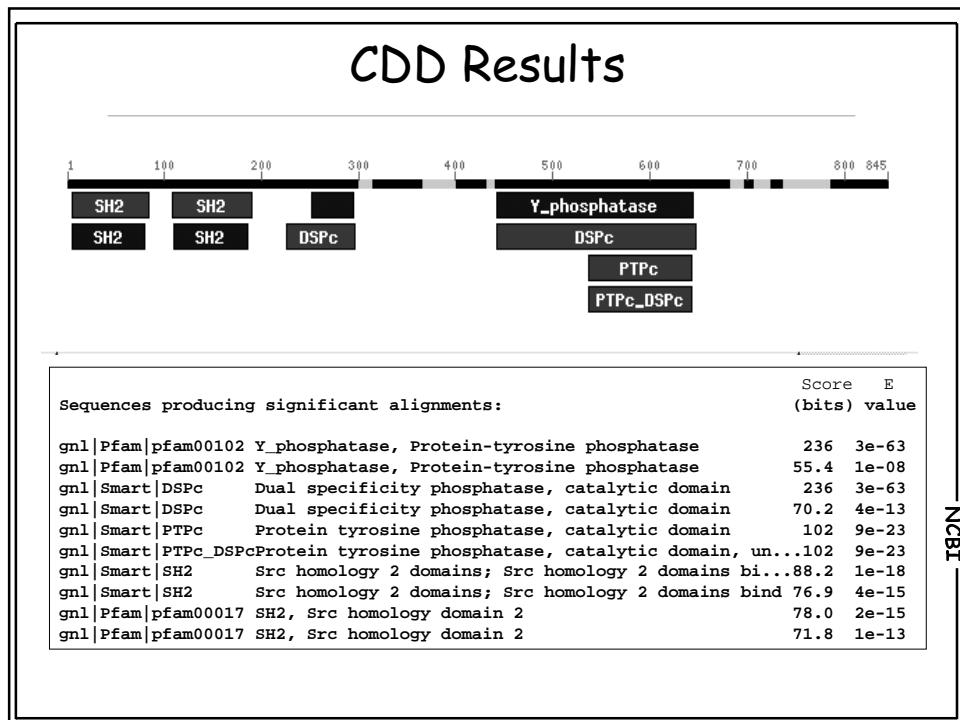
Search the Conserved Domain Database with Reverse Position Specific BLAST

**Reference:** Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

Search Database:    
Enter query as Protein  Sequence in FASTA format

```
>gi|7290263|gb|AAF45724.1| CG3954 gene product [alt 2] [Drosophila melanogaster]
MSRRWFHPTISGIEAEKLLQEQQGFDGSFLARLSSSNPGAFTLSVRRGNEVTHIKIQNNGDF
FDLYGGEEKFATLPELVQYYMENGELEKEKNQAAIELKQPLICAEPPTTERWFHGNLSGKEAEKL
ILERGKNGSFVLVRESQSKPGDFVLSVRTDDKVTHVMIRWQDKKYDVGGGESFGTLSELIDHY
KRNPMPVETCGTVVHLRQPFPNATRITAAGINARVEQLVKGGFWEEFESLQQDSRDTFSRNEGY
EQENRLKNRYRNILPYDHTRVKLLDVEHSVAGAEYINANYIRLPTDGDLYNMSSSESLSNNS
VPSCPACTAAQTQRNCSNCQLQNKTCVQCAVKSAILPYNSNCATCSRKSDSLSKHKRSSESSAS
```

NCBI



## Options for Advanced Blasting: nucleotide

**Options** for advanced blasting

Limit by entrez query  AND (none)

**Example Entrez Queries**

- nucleotide all[Filter] NOT mammalia[Organism]
- green plants[Organism]
- biomol mRNA[Properties]
- biomol genomic[Properties]

**Other**

**OtherAdvanced**

- W 7 word size to 7
- e 10000 expect value
- v 2000 descriptions
- b 2000 alignments

only  Mask lower

## Options for Advanced Blasting: protein

Options for advanced blasting

Limit by entrez query  AND  (none)

Composition-based statistics

### Example Entrez Queries

proteins all[Filter] NOT mammalia[Organism]  
green plants[Organism]  
srcdb refseq[Properties]

### OtherAdvanced

-W 2 word size to 2  
-e 10000 expect value  
-v 2000 descriptions  
-b 2000 alignments

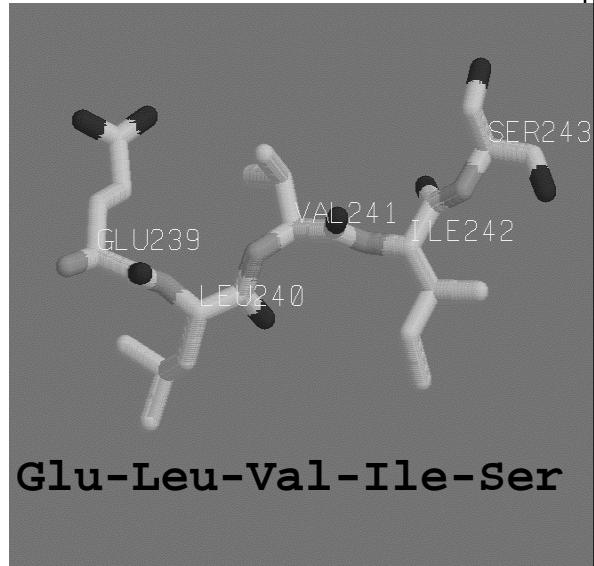
Organism search pull-down



NCBI

PHI pattern

## The Elvis Problem: Short Sequences BLAST



## Finding Hits with Short Sequences

Composition-based statistics

Choose filter  Low complexity  Mask for lookup table only  Mask lower case

Expect

Word Size

Matrix  Gap Costs Existence: 9 Extension: 1

**245 Elvises**  
Sequences producing sign

Query: 1 ELVIS 5	(bits) Value
<b>ELVIS</b>	ply... 18 7266
<b>Sbjct: 55 ELVIS 59</b>	hil... 18 7266
	ote... 18 7266
	TRANSSCRIPTION-REPA... 18 7266
	probable 8-oxoguan... 18 7266
	(AE003698) CG7472 gene product [D... 18 7266
	(XM_080244) CG1918 [Drosophila... 18 7266
	colonization factor antigen I, CFA... 18 7266
	(AL161931) bA1021019.1 (zinc fi... 18 7266
	(L76577) immunoglobulin light ch... 18 7266
	(AB029894) p3vc [rice grassy stu... 18 7266
	(L24774) delta3, delta2-enoyl-CoA ... 18 7266
	(NC_003272) cobalt transport p... 18 7266

NCBI

## Monkeys Typing

```
>dbj|BAB12211.1| (AB032549) polyketide synthase and peptide synthetase [Microcystis aeruginosa]
Length = 3487

Score = 22.2 bits (45), Expect = 294
>gb|AAF63753.1|AF242291_1 (AF242291) nuclear protein EAST [Drosophila melanogaster]
Length = 2362

Score = 23.5 bits (48), Expect = 120
Identities = 8/12 (66%), Positives = 9/12 (74%), Gaps = 1/12 (8%)
>pdb|1A59| Cold-Active Citrate Synthase
Length = 378

Score = 23.5 bits (48), Expect = 120
Identities = 8/12 (66%), Positives = 9/12 (74%), Gaps = 1/12 (8%)
>gb|AAE00000.1|AAE00000_1 (AAE00000) putative protein [Unknown]
Length = 1000

Query: 1 ELVISPRESLEY 12
EL I PRE L+Y
Sbjct: 145 EL-IEPREDLDY 155
```

NCBI

## Large Sequences

>193,787 bases

GAATTCAAGTGTCTTATTCAGGTAATACGGCAATTGTACTAGTTGGAAATGAAATATCAA  
GGCTTCAGTTAACGAGTCATTATGTATTAAAGCTGTCTGGCTGTTAAAAGA

Other advanced options:

-e 1e-34 -W 30

TATCTTTGCAAACCATTAAAGAGTTTCTGACTGACAAAAATAAGTAGGCTGGAAGTTAAAAGAA  
ATGCTCGCTTCTGACAGGATCTAGCTTAGATGTTCTATACGGTTGATAATGCCAGCCCT  
AAGCTATTCAATCTTTTCTTTCTTTCTTTAAAACCTTCTGACTAGATCAAGACA

**ERROR: Bl: Megablast nit exceeded**

ATAAAAAATTTTGCTGGTAAATATAAAAGAAATCTGAATAATTAAAAGTAATAAGGGAAACAU  
ATTCAAGAATGATCTCATTCTCAAACCGTTGTTCAAAGATATGTTAAGACATTAAAGGACTCC  
AGTGATGAATTAATTCTTAAACATGAGCAATATTAAATGATAGAAATTATAACAGGAA  
ATGTTGTAATTCTATGGTAGTTCTAGTGGAAAAGTATAGCATCCCTGGCTGGCACAGTG  
ATACCTGTAATCCAGTGTCTGGGAGGAGCTGGAGGATAGCTCGAGGCCAGAAGTTGAGA  
CCTGGCAGCATAGCAAGACCGTGTCTTACAAACGATTTTATATATGTGTATATATTCAAT  
TATAAATATGTAATAATTATATAAAATATGTAATAATTATACATAAAA

NCBI

## MegaBlast

> 4788 gnl|UG|Os#S4788 96BS0324 Oryza sativa cDNA /clone=96BS0324  
GNAATTGAAATACGACTCACTATAGGGCAATTGGTACCGGGCCCCCTCGAGTTTT  
TTTAACTTTTGTAACTGAAATCTCGGATACTAAAGTTATAAAAGGGAAACTA  
GCTAACATTCTCCATACATCATCCAGTACCATAGTAAGGCTGCTGCTAGTTGCATAGCC  
CGATAAGAGTCTCACACAAGCACAGAAGGTTAAGAATGGGAAAGACGAAAAATCACA  
GAGAACAGTAAACATTCTAGAGCTAGCTAACGCTTTCTGCTCATCCCCATTCTG  
CCTTTGCAATGCCAACGGCTGCTCGCCCTCGCAAGCTCCCTCAGTTACTCCAGG  
ATCATGTTCTCGGGGGCACCCAAACGGCAGGGCTCATGATGTTCAAGCATTGCC  
CTCTGCTTGCAGGGCTTACCACTGGTGTGCAAAGAACAGGGGTGCCTTGGT  
AAGTACTCAGGATGGTANCCACTGGTGGAAAGATGGACTCCCNCCCCCGNTTCACT  
GATCTGGTCTAACTCGGAAGGACAAATCAAAGATCGGGCGGAAGGATNATCNCA  
GGTTTNTCAAAAATGGCTACCCANAAATTGAGTNNTCTATGCCTGNTCCAAAATT  
> 70988 gnl|UG|Os#S70988 H061C07 Oryza sativa cDNA /clone=H061C07  
GGCTACCATCTGAGCTACCTCACCAAGGCACCCCTGTTCTCTGCCACACAGATGGT  
GAAGGCCCAGGGCAAGCAGAGGGCAATCTGGAGAACATCATGAGGGCCTGGCTGGC  
TGCCACGGGAGAACACATGATCTGGAGTACAAGTGAGGGAGGCTGCGAGGACAG  
CACCCGGGGCATGGCAAGAGGGCAATGGGATGGCAAGAACAGCTGATTAGCTAG  
CTCTAAGAATTTAGCTCTGGTATTTCTGGCTTTCCATTCTAACCTTCT  
GTGCTTGTGAGACTCTTATGGGCTATGCAACTCAGCAGCTTACTATGGTACTGG  
ATGATTTATGGAGAATGGTAGCTAGTCCCTCTATTATAACTTATTAGTATGGAGA  
TTTCAAGTCAAAA  
> 69736 gnl|UG|Os#S69736 H030C04 Oryza sativa cDNA /clone=H030C04  
AGCTTACCTCACCAAGGCACCCCTGTTCTCCATGCCACACAGATGGTGAACGCCCTGGCA  
AAGCAGAGGGCAATGGCTGGAGAACATCATGAGGGCCTGGCTGGCTGGCCCCGAGAAC  
ACATGATCTGGAGTACAAGTGAGGGAGGCTGGAGGGAGACCCGGCATGTGGCAA  
AGAGGGCAATGGGATGGAGAACAAAGACGTGATTAGCTAGCTAAAGATTGTTTAGC  
TTCTCTGTGATTTCTGGCTTCCATTCTTAACCTCTGTGCTTGTGAGACTCT  
TATCGGGCTATGCAACTAGCAGCAGCTTACTATGGTACTGGATGATGTTATGGAGAATG  
TTAGCTAGTCCCTCTATTATAACTTATTAGTATGGAGATTCAAGTCAAAAAAAA  
AAAAAA  
> 58169 gnl|UG|Os#S58169 AU063597 Oryza sativa cDNA /clone=C63051\_1A  
TACCTCTGGATNCGGCCGGCCGGCCGGCACTGGGATGCCAACAGAACCTGATTAGCTA

NCBI

## BLAST: standalone, clients, databases

Current directory is /blast

Up to higher level directory

```
└── README  
└── blasturl/  
└── db/  
└── documents/  
└── executable/  
└── fmerge/  
└── matrices/  
└── network/  
└── old/  
└── server/  
└── temp/
```

```
ftp> open ftp.ncbi.nih.gov
```

```
.
```

```
.
```

```
ftp> cd blast
```

ftp://ftp.ncbi.nih.gov/blast/

NCBI

## BLAST Batch Client

```
C:\Netblast>blastcl3 -i input.seq -d nr -p blastn -o outfile
```

```
National Center for Biotechnology Information (NCBI)
```

```
welcome to the blast network service.  
BLASTN 2.2.2 [Dec-14-2001]
```

```
Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,  
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),  
"Gapped BLAST and PSI-BLAST: a new generation of protein database search  
programs", Nucleic Acids Res. 25:3389-3402.
```

```
Query= gi|19343352|gb|AAF15280.2|AF192502_1 aryl hydrocarbon receptor  
[Gallus gallus]  
(535 letters)
```

```
Database: All GenBank+EMBL+DDBJ+PDB sequences (but no EST, STS, GSS,  
or phase 0, 1 or 2 HTGS sequences)  
1,214,651 sequences; 1,071,392,519 total letters
```

NCBI

# Genomic BLAST Pages

## Genomic BLAST pages [?](#)

- [Human Genome](#)
- [Mouse Genome](#)
- [Rat Genome](#)
- [Fugu rubripes](#)
- [Zebrafish Genome](#)
- [Anopheles gambiae](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Other eukaryotes](#)
- [Microbial Genomes](#)

N  
NCBI

## Microbial Genomes BLAST

### BLAST with microbial genomes (126 bacterial)

Currently available for BLAST searches are sequences from which GenBank has received a completed genomic sequence

for which only a partial genome is available.

See [About the Database]

P - indicates the ability to search  
- unfinished genomic sequence

Enter your query sequence:

```
>APE0122
MVGVFGRRLSRHVVVKRW
VAELEDLTMDTEIVEMAK
TFKKLKPRAVEEEARKEG
```

Select type of query and

Query:  Protein  DNA

You may change BLAST

Expect: 10

1

### Pasteurellaceae

- Actinobacillus actinomycetemcomitans*
- Haemophilus ducreyi*

P  *Haemophilus influenzae Rd*

P  *Pasteurella multocida*

### Pseudomonas

- P  *Pseudomonas aeruginosa*
- P  *Pseudomonas fluorescens*
- Pseudomonas putida KT2440*
- Pseudomonas putida PRS1*
- Pseudomonas syringae*

### Salmonella

- Salmonella enterica* subsp. *enterica* serovar Dublin
- Salmonella enteritidis*
- Salmonella paratyphi*
- Salmonella typhi*
- Salmonella typhimurium LT2*

### Spirochaetales

- P  *Borrelia burgdorferi*
- Treponema denticola*
- P  *Treponema pallidum*

### Thermotogales

- P  *Thermotoga maritima*
- P  *Thermus/Deinococcus group*

### Deinococcus radiodurans

## Hits to Unfinished Genome

Sequences producing significant hits in query:

```
gnl|TIGR_61435|6428 Dehalococcobacter
ref|NC_00918.1| Aquifex aeolicus
gnl|DOE_20211|2351476 fasta.scre
gnl|DOE_2108|2351478 fasta.scre
gnl|TIGR_881|11531 Desulfobacter
ref|NC_002677.1| Mycobacterium
gnl|TIGR_1772|3270 Mycobacterium
emb|AL137778.1|SCL2 Streptomyces
gnl|CBCCMH_4770|parabactin.fasta.sc
gnl|TIGR_35554|2947 Geobacter
gnl|TIGR_1764|3294 Mycobacterium
ref|NC_002755.1| Mycobacterium
gnl|Sanger_1765|mbovis Contig20
gnl|TIGR_164513|mtub210_338 Mycobacterium
ref|NC_000962.1| Mycobacterium
emb|AX127148.1|AX127148 Sequen...
gnl|Sanger_17171|Corynebacterium
ref|NC_002163.1| Campylobacter
gnl|TIGR_10971|3499 Chlorobium
ref|NC_000915.1| Helicobacter
ref|NC_000921.1| Helicobacter
gnl|SANGER_36826|Contig101 Clostridium
gnl|DOE_63737|Contig554 Nostoc
ref|NC_003030.1| Clostridium
ref|NC_002570.1| Bacillus halophilus
gnl|SANGER_36826|Contig95 Clostridium
gnl|TIGR_13921|6277 Bacillus anatum
gnl|TIGR_13921|gbal18 Bacillus anatum
ref|NC_003296.1| Ralstonia solanacearum
gnl|TIGR_164513|mtub210_bmt782 Mycobacterium
gnl|TIGR_164513|mtub210_449 Mycobacterium
gnl|Sanger_1765|mbovis Contig27
```

**Contributing Genome Centers**

- Entrez Genomes
- ASTRA
- SubtiList
- BSNR
- ChGP
- Columbia
- Diversa
- Genome Therapeutics
- Gencore
- Gencore Genoscope
- GTC
- Heidelberg
- KDRI
- NITE
- OU-ACGT

**Thermobifida fusca, unfinished sequence**

Accession: NC\_002721  
Number of Contigs: 64  
Total Bases Sequenced: 3656222 bp  
Last updated: Mar 12, 2002.

**Contributor: DOE-JGI**

**Organism: Thermobifida fusca**  
Genetic Code: 11  
Lineage: Bacteria; Firmicutes; Actinobacteria; Actinobacteridae; Actinomycetales; Streptosporangineae; Nocardiopsaceae; Thermobifida.

**COMMENT** JGI Data Release Policy: The Joint Genome Institute releases data as early and often as updates are feasible. All data is available to our scientific colleagues for the purpose of searching for genes of interest for downstream biological characterization. Please recognize that these data are preliminary, and therefore may contain errors, and are subject to change (particularly with respect to contiguity) in subsequent releases. The JGI asks that you acknowledge the source of information obtained from this site in any publication benefiting from these data by including the following sentence in both the Materials and Methods and Acknowledgement sections: 'Preliminary sequence data was obtained from The DOE Joint Genome Institute (JGI)' at

## BLAST with At Genome

Arabidopsis thaliana

BLAST Home Page

BLAST overview

BLAST FAQs

BLAST news

BLAST manual

**BLAST with Arabidopsis thaliana genome**

Compare your query sequence to genomic DNA sequence, mRNAs or protein products.

Database: clone ▾ Program: tblastn ▾ MegaBlast: □

Begin Search

Enter an accession, gi, or a sequence in FASTA format:

```
>gi|13878337|sp|Q13315|ATM_HUMAN Serine-protein kinase ATM (Atm)
MSLVLDLICCRQLEHDRTERRKKEVEKFKKLIRDPETIKHLDRHSDSKQGKYLNWDAV
ETECLRIAKPWNVASTQASPKWMEISSLVKYFICKANRRAAPRLKCQEELLNYIMDTVKD
CSNILLKDILSVRKTKMCEISQQQMELFSPVYFRDYLKPSQDVHRVIVARIIIHAVTKGCCS
DFFSKAIQCQCARQEKSSSGLNHHLAAATIPLKTLLAVNFRIRVCGLGDEILFTLLYIWTQHR
ELFQLQIYIHHPGAKTQEKGAYESTKWRSLILYMLYDLLVNEISHIGSRGKYSQGFRNIA
```

Optional parameters

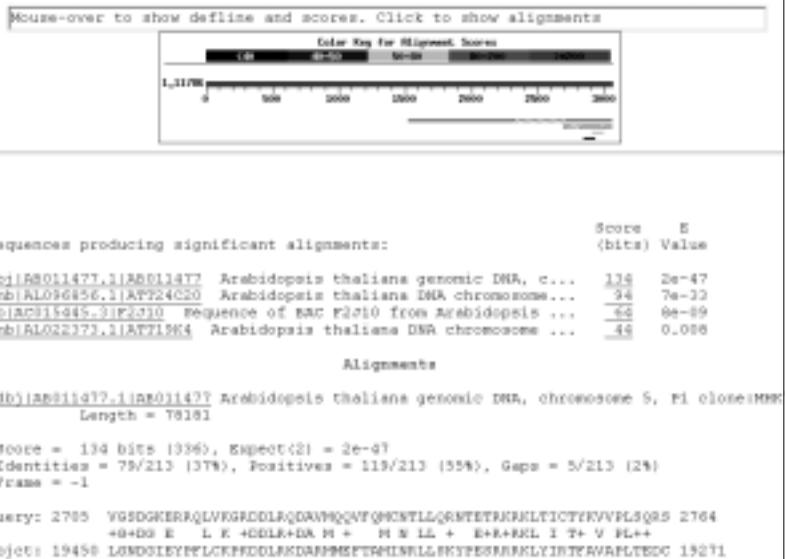
Expect	Filter	Descriptions	Alignments
0.01	default	100	100

Advanced options:

Begin Search Clear Input

## Hits to At Genome

### Distribution of 10 Blast Hits on the Query Sequence



NCBI

## Genomic Context of BLAST Hits

### Master Map: Clone

### Maps & Options

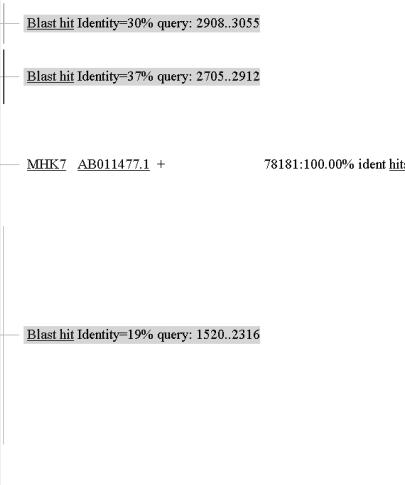
Total Clones On Chromosome: **416**

Region Displayed: **16,057K-16,063K**

**b** Download/View Sequence/Evidence

Clones Labeled: **4 Total Clones in Region: 4**

Gene Clone clone accession orientation align.status



NCBI

**The Rice Genome**

Oryza sativa

---

[BLAST Home Page](#)

[BLAST overview](#)

[BLAST FAQs](#)

[BLAST news](#)

[BLAST manual](#)

**BLAST against Oryza sativa ssp. indica WGS contigs**

In that none of the contigs have been mapped to a rice genome, we are unable to provide a display of the location of any BLAST hits in the rice genome.

Database:  contig    Program:  blastn    MegaBlast    Begin Search

Enter an accession, gi, or a sequence in FASTA format:  

```
>gi|13878337|sp|Q13315|ATM_HUMAN Serine-protein kinase ATM (MSLVNLNDLICCRQLEDRATERKKEVEKFRLIRDPEТИKHLDRHSDSKQGKYLNNDAVETECLRIAKPNVSASTQ&SROKKMMEISSLVVKYFIRCANR&PRLKCQELLNYIMDTVKDCSNSNLLKDILSVRKYWCIESQQQULELFSVYFRLYLKPSQDVHRVLVARIIHAVTKGCSDFFSKAIQCARQEKSSSGLNHILAALTIFLKLTLAVNFRIRVCCELGDEILPTLLYIWTQHRELFQLQIYIHHPKGAKTQEKGAYESTKWRSLVNLYDLLVNEISHIGSRGKYSSGFRNIA
```

Optional parameters  
 Expect    Filter    Descriptions    Alignments  
 0.01    default    100    100  
 Advanced options:

NCBI

**Shotgun Contigs**

RID: 1019619759-0

Query= gi|1387833  
ATM (ataxia telangiectasia, telomeric locus) (3056 bp)

Database: Oryza sativa (103,040 contigs)

If you have any problems with this page, please refer to the [FAQ](#).

[Taxonomy reports](#)

[Mouse-over to show sequence](#)

[1\\_13834](#)

**LOCUS** AAAA01003283 20559 bp DNA linear PLN 04-APR-2002  
**DEFINITION** Oryza sativa (indica cultivar-group), whole genome shotgun sequence.  
**ACCESSION** AAAA01003283  
**VERSION** AAAA01003283.1 GI:19927592  
**KEYWORDS** .  
**SOURCE** Oryza sativa (indica cultivar-group).  
**ORGANISM** Oryza sativa (indica cultivar-group)  
**Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta; Spermatophyta; Magnoliophyta; Liliopsida; Poales; Poaceae; Ehrhartioideae; Oryzeae; Oryza.**

**REFERENCE** 1 (bases 1 to 20559)  
**AUTHORS** Yu,J., Hu,S., Wang,J., Wong,G.K.-S., Li,S., Liu,B., Deng,Y., Dai,L., Zhou,Y., Zhang,X., Cao,M., Liu,J., Sun,J., Tang,J., Chen,Y., Huang,X., Lin,W., Ye,C., Tong,W., Cong,L., Geng,J., Han,Y., Li,L., Li,W., Hu,G., Huang,X., Li,W., Li,J., Liu,Z., Li,L., Liu,J., Qi,Q., Liu,J., Li,L., Li,T., Wang,X., Lu,H., Wu,T., Zhu,M., Ni,P., Han,H., Dong,W., Ren,X., Feng,X., Cui,P., Li,X., Wang,H., Xu,X., Zhai,W., Xu,Z., Zhang,J., He,S., Zhang,J., Xu,J., Zhang,K., Zheng,X., Dong,J., Zeng,W., Tao,L., Ye,J., Tan,J., Ren,X., Chen,X., He,J., Liu,D., Tian,W., Tian,C., Xia,H., Bao,Q., Li,G., Gao,H., Cao,T., Wang,J., Zhao,W., Li,P., Chen,W., Wang,X., Zhang,Y., Hu,J., Wang,J., Liu,S., Yang,J., Zhang,G., Xiong,Y., Li,Z., Mao,L., Zhou,C., Zhu,Z., Chen,R., Hao,B., Zheng,W., Chen,S., Guo,W., Li,G., Liu,S., Tao,M., Wang,J., Zhu,L., Yuan,L. and Yang,H.

**TITLE** Draft Sequence of the Rice Genome (Oryza sativa L. ssp. indica)  
**JOURNAL** Science 296, 79-92 (2002)

Sequences produced:

gb AAAA01003283.11	Oryza sativa (indica cultivar-group) sca...	139	4e-31
gb AAAA010032561.11	Oryza sativa (indica cultivar-group) sca...	124	1e-26
gb AAAA010095014.11	Oryza sativa (indica cultivar-group) sca...	82	4e-21
gb AAAA010707688.11	Oryza sativa (indica cultivar-group) sca...	102	4e-20
gb AAAA01001778.11	Oryza sativa (indica cultivar-group) sca...	53	5e-05

NCBI

## Service Addresses

• **General Help**

info@ncbi.nlm.nih.gov

• **Questions about BLAST**

blast-help@ncbi.nlm.nih.gov

**E-mail Servers**

**BLAST Server**

blast@ncbi.nlm.nih.gov

NCBI